

1 **The machine learning algorithm identified COL7A1 as a diagnostic marker for**  
2 **LUSC and HNSC**

3 Chenyu Wang<sup>1</sup>; Yongxin Ma<sup>2</sup>; Jiaojiao Qi<sup>2,3\*</sup>; Xianglai Jiang<sup>4\*</sup>

4 <sup>1</sup> Xiangya Hospital of Central South University, Changsha, 410008, China

5 <sup>2</sup> Ningxia Medical University, Yinchuan, 750004, China

6 <sup>3</sup> General Hospital of Ningxia Medical University, 750004, China

7 <sup>4</sup> School of basic medical sciences and life sciences, Hainan Medical University,  
8 Haikou, 571199, China

9 **\* Correspondence:**

10 Jiaojiao Qi

11 [qjjdemyx@163.com](mailto:qjjdemyx@163.com)

12 Xianglai Jiang

13 [jiangxl1997@126.com](mailto:jiangxl1997@126.com)

14 **Abstract**

15 Squamous cell carcinomas (SCCs) comes from different parts, but there may be similar  
16 tumorigenic signaling pathways and metabolism, and different squamous cell  
17 carcinoma has a similar mutation landscape and squamous differentiation expression.  
18 Studying the expression profile of common SCCs is helpful to find biomarkers with  
19 diagnostic and prognostic significance for a variety of squamous cell carcinoma. Lung  
20 squamous cell carcinoma (LUSC), head and neck squamous cell carcinoma (HNSC),  
21 and ‘squamous cell cancer’ in esophageal carcinoma (ESCA) and cervical squamous  
22 cell carcinoma and endocervical adenocarcinoma (CESC) in The Cancer Genome Atlas  
23 (TCGA) database were used as training sets. The relevant data sets in the Gene  
24 Expression Omnibus (GEO) database were selected as validation sets. Machine  
25 learning algorithms were used to screen out factors with high accuracy in the diagnosis  
26 of SCCs as core genes, and explore their effects on patient prognosis and  
27 immunotherapy. COL7A1 (Collagen Type VII Alpha 1 Chain) has high accuracy in the  
28 diagnosis of LUSC and HCSC, whether in the training set (LUSC \_AUC: 0.987; HNSC  
29 \_AUC: 0.933) or validation set (LUSC \_AUC: 1.000; HNSC \_AUC: 0.967). Moreover,  
30 the expression of COL7A1 was significantly correlated with shorter OS and DSS in  
31 HNSC and LUSC patients, and was also significantly negatively correlated with IPS in  
32 LUSC patients treated with CTLA4 (-) PD1 (+), CTLA4 (+) PD1 (-) and CTLA4 (+)  
33 PD1 (+). COL7A1 has the potential to be used as a diagnostic and prognostic marker  
34 for LUSC and HNSC and to predict the efficacy of LUSC immunotherapy.

**NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

35 **Keywords:** *machine learning, squamous cell carcinoma, diagnostic biomarker,*  
36 *COL7A1*

### 37 **Introduction**

38 Squamous cell carcinomas (SCCs) occur in epithelial cells and squamous cells.  
39 Common lesions include skin, esophagus, cervix, vagina, bronchus, etc<sup>1-6</sup>. The known  
40 causes of squamous cell carcinoma are diverse, including gene mutations, age,  
41 ultraviolet and human papillomavirus infection<sup>3, 7</sup>. SCCs with high incidence mainly  
42 include esophageal squamous cell carcinoma (ESCC), head and neck squamous cell  
43 carcinoma (HNSC), lung squamous cell carcinoma (LUSC) and cervical squamous cell  
44 carcinoma (CSCC). Lung squamous cell carcinoma is derived from bronchial  
45 epithelium and belongs to non-small cell lung cancer (NSCLC), and the 5-year survival  
46 rate of patients with LUSC is less than 20 %<sup>8</sup>. The diagnosis of LUSC and lung  
47 adenocarcinoma (LUAD) can be identified by detecting whether the  
48 immunohistochemical staining results of P63, P40 and CK5 / 6 are positive<sup>9</sup>. HNSC is  
49 the sixth most common cause of cancer death in the world. Due to the high recurrence  
50 and metastasis rate of advanced HNSC, the prognosis of patients with advanced HNSC  
51 is still poor even with various treatments such as surgery, radiotherapy and  
52 chemotherapy<sup>10</sup>. Esophageal squamous cell carcinoma ESCC is the main subtype of  
53 esophageal cancer. ESCC is highly heterogeneous, and the 5-year survival rate is about  
54 20 %<sup>11, 12</sup>. CESC is the second most common gynecological cancer; there is still a lack  
55 of good biomarkers to detect early CESC, and CESC is often found in the advance  
56 stage<sup>13</sup>. SCCs come from different parts, but may have similar tumorigenic signaling  
57 pathways and metabolism, and different SCCs have similar mutation landscapes (such  
58 as TP53, SOX2 and TP63) and squamous differentiation expression<sup>14-18</sup>. Studying the  
59 expression profile of common SCCs is conducive to the discovery of biomarkers with  
60 diagnostic and prognostic significance for a variety of squamous cells, to detect SCCs  
61 early to improve the therapeutic effect of SCCs patients.

62 Bioinformatics technology provides a new method for the screening of tumor markers.  
63 The analysis of high-throughput sequencing data using bioinformatics technology is  
64 more conducive to finding genes with diagnostic value. The machine learning algorithm  
65 LASSO regression and SVM-RFE were used to identify the core genes with diagnostic  
66 potential. The Receiver Operating Characteristic Curve (ROC curve) was drawn and  
67 the area under the ROC curve (AUC) was calculated. The core genes with the most  
68 diagnostic value were selected as the diagnostic markers of squamous cell carcinoma,

69 and the possible role of core genes in the development of specific cancers was  
70 analyzed<sup>19-21</sup>.

## 71 **Method**

### 72 **Collect and process datasets**

73 The RNA-seq and corresponding clinical data of LUSC, HNSC, and SCCs samples in  
74 ESCA and CESC in the TCGA database were used as training sets. Similarly, datasets  
75 from GEO (GSE9850, GSE19188, GSE23400, GSE30784) were collected and the sva  
76 R package was used to perform batch correction for four datasets<sup>22-25</sup>. The Limma R  
77 software package was used to calculate the differential genes between squamous cell  
78 carcinoma tissues and adjacent tissues. Differential genes with  $|\logFC| > 2$  and p value  
79  $< 0.05$  were selected as validation sets.

### 80 **Machine learning algorithms screen core genes**

81 LASSO regression can use the constructed penalty function to calculate the optimal  
82 model and prevent the model from overfitting. The genes formed by the optimal model  
83 analyzed by lasso regression were used as the core genes for the diagnosis of squamous  
84 cell carcinoma. SVM-RFE is a machine learning algorithm based on Embedded method,  
85 which can complete the sorting of feature genes while screening useful feature genes<sup>19-</sup>  
86 <sup>21</sup>. The LASSO model was used to screen the core genes with diagnostic value for  
87 squamous cell carcinoma. On this basis, the SVM-RFE algorithm was used to calculate  
88 the final gene set with diagnostic value. The ROC R software package was used to draw  
89 the ROC curve of each gene expression in the core gene set and calculate the area under  
90 the ROC curve. Genes with high accuracy in both training set and validation set were  
91 selected as diagnostic markers for squamous cell carcinoma. The AUC of the expression  
92 of diagnostic markers in HNSC, LUSC, CSCC and ESCC was calculated respectively,  
93 and cancers with high diagnostic rates were selected to further analyze the possible  
94 biological functions of diagnostic markers in these cancers.

### 95 **Clinical features and prognosis analysis**

96 The expression of core genes in squamous cell carcinoma was selected, and the  
97 correlation between the expression of core genes and the patient 's age, gender, stage, T  
98 staging, N staging, M staging was completed. The data of OS, DSS, DFI and PFI in  
99 patients with squamous cell carcinoma were collected and divided into two groups  
100 according to the expression of core genes. The correlation between the expression of  
101 SCCs core genes and OS, DSS, DFI and PFI was completed.

### 102 **Methylation analysis**

103 Methylation is a kind of epigenetic phenomenon and the main form of epigenetics<sup>26</sup>.  
104 CpG islands are generated by promoter methylation to silence genes, thereby regulating  
105 gene expression. Methylation does not change DNA sequences, but can change traits  
106 and pass on to the next generation<sup>27</sup>. The promoter methylation data of core genes in  
107 UALCAN website were collected to study the effect of core gene methylation on OS  
108 in patients with squamous cell carcinoma.

### 109 **Enrichment analysis**

110 Gene ontology (GO) analysis involves biological processes, molecular functions and  
111 cellular components. Kyoto Encyclopedia of Genes and Genomes (KEGG) analysis, as  
112 a widely used database, includes information on currently known signaling pathways.  
113 The 'c5.go.v7.4.symbols.gmt' file and the 'c2.cp.kegg.v7.4.symbols.gmt' file from the  
114 Molecular Signatures Database (MSigDB) were downloaded and used to perform  
115 GSEA-based GO enrichment analysis and KEGG enrichment analysis on differential  
116 genes in high and low expression samples of core genes.

### 117 **Analysis of tumor microenvironment and immune cell infiltration**

118 The tumor microenvironment includes the structure, function and metabolism of the  
119 tissue where the tumor is located. It is also related to the internal environment of the  
120 tumor cells themselves. The tumor microenvironment regulates the occurrence,  
121 development and metastasis of the tumor. As a complex and dynamically regulated  
122 system, tumor microenvironment is composed of tumor cells, immune cells and  
123 supporting cells. The Estimate R software package was used to calculate the correlation  
124 between the expression of core genes and StromalScore, ImmuneScore and  
125 ESTIMATEScore in specific squamous cell carcinomas using wilcox test. The TMIER  
126 method was used to calculate the correlation between core gene expression and  
127 infiltration of B cell, T cell CD4, T cell CD8, Neutrophil, Macrophage and DC.

### 128 **Analysis of Immunotherapy**

129 As an indicator of tumor immunogenicity, immunophenoscore (IPS) can predict the  
130 effect of immunotherapy in cancer patients<sup>28</sup>. Immunotherapy data were obtained from  
131 the TCIA database, and the correlation between the expression of core genes and IPS  
132 treated with CTLA4 (-) PD1 (-), CTLA4 (-) PD1 (+), CTLA4 (+) PD1 (-) and CTLA4  
133 (+)

## 134 **Results**

### 135 **Determination of core genes**

136 The limma package was used to screen the differential genes (227 in total). The results

137 of lasso regression analysis showed that a total of 37 genes were identified (Fig.1.A).  
138 The svm-rfe algorithm showed that a total of ten genes (SPP1, VEGFD, IL33,  
139 GPIHBP1, KIF18B, CRNN, HJURP, COL7A1, TMEM132A, CHI3L2) were identified  
140 as core genes (Fig.1.B). The AUC results of each gene in SCCs showed that  
141 TMEM132A (AUC: 0.944), KIF18B (AUC: 0.975), HJURP (AUC: 0.976) and  
142 COL7A1 (AUC: 0.929) had a higher diagnostic rate for squamous cell carcinoma in the  
143 training set (TCGA dataset) (Fig.2.A-J). The analysis results in the validation set (GEO  
144 dataset) showed that only COL7A1 (AUC: 0.907) had a high diagnostic rate for SCCs,  
145 so COL7A1 was selected as the core gene for the diagnosis of squamous cell carcinoma  
146 (Fig.3.A-D). The analysis of COL7A1 in different squamous cell carcinomas showed  
147 that the AUC of COL7A1 in CESC was 0.968 (Fig.4.A), the AUC in ESCA was 0.537  
148 (Fig.4.B), the AUC in HNSC was 0.933 (Fig.4.C) and the AUC in LUSC was 0.987  
149 (Fig.4.D) in the TCGA dataset. The AUC of COL7A1 in CESC was 0.648 (Fig.4.E), in  
150 ESCA was 0.951 (Fig.4.F), in HNSC was 0.967 (Fig.4.G) and in LUSC was 1.000  
151 (Fig.4.H) in GEO dataset. COL7A1 was selected as a diagnostic biomarker for HNSC  
152 and LUSC for further analysis.

### 153 **Clinical features and prognosis analysis**

154 The results of clinical characteristics analysis showed that there was no significant  
155 correlation between the expression of COL7A1 and age, gender, stage, T staging, N  
156 staging and M staging in LUSC patients and HNSC patients (Fig.5.A-F). The results of  
157 KM prognostic analysis showed that the high expression of COL7A1 was significantly  
158 positively correlated with shorter OS (Fig.6.A) and DSS (Fig.6.C) in LUSC patients  
159 and shorter OS (Fig.6.E) and DSS (Fig.6.G) in HNSC patients. There was no significant  
160 difference in DFI (Fig.6.B) and PFI (Fig.6.D) between LUSC patients with high and  
161 low expression of COL7A1, and there was no significant difference in DFI (Fig.6.F)  
162 and PFI (Fig.6.H) between HNSC patients with high and low expression of COL7A1.

### 163 **Promoter methylation level**

164 The methylation level of COL7A1 in LUSC and HNSC obtained from UCSC xena  
165 website showed that COL7A1 had a lower Promoter methylation level in LUSC tumor  
166 (Fig.7.A-B).

### 167 **Enrichment analysis**

168 Go enrichment analysis showed that the expression of COL7A1 was positively  
169 correlated with the activation of artery development, detection of chemical stimulus,  
170 sensory perception of smell, odorant binding and olfactory receptor activity in LUSC

171 (Fig.8.A). And COL7A1 was positively correlated with the activation of etection of  
172 chemical stimulus, muscle cell proliferation, sensory perception of smell, transforming  
173 growth factor beta receptor signaling pathway and olfactory receptor activity in HNSC  
174 (Fig.8.C). KEGG enrichment analysis showed that the expression of COL7A1 was  
175 positively correlated with the activation of olfactory transduction, and negatively  
176 correlated with the activation of folate biosynthesis in LUSC (Fig.8.B). The expression  
177 of COL7A1 was positively correlated with the activation of olfactory transduction, and  
178 negatively correlated with the activation of O glycan biosynthesis, riboflavin  
179 metabolism and starch and sucrose metabolism in HNSC (Fig.8.D).

### 180 **Tumor microenvironment and immune cell infiltration**

181 The results of tumor microenvironment analysis showed that the StromalScore,  
182 ImmuneScore and ESTIMATEScore of LUSC samples in the low expression group of  
183 COL7A1 were significantly higher than those in the high expression group of COL7A1  
184 (Fig.9.A). There was no significant difference in StromalScore, ImmuneScore and  
185 ESTIMATEScore between HNSC samples of COL7A1 low expression group and  
186 HNSC samples of COL7A1 high expression group (Fig.9.B). The results of immune  
187 cell infiltration analysis showed that the expression of COL7A1 was significantly  
188 negatively correlated with the infiltration of B cell, T cell CD8, Macrophage and DC in  
189 LUSC, and the expression of COL7A1 was significantly negatively correlated with the  
190 infiltration of B cell, T cell CD4, Macrophage and DC in HNSC (Fig.9.C).

### 191 **Immunotherapy**

192 Immunotherapy analysis from the TCIA database showed that the expression of  
193 COL7A1 was not significantly correlated with IPS treated with CTLA4 (-) PD1 (-)  
194 (Fig.10.A), and was significantly negatively correlated with IPS treated with CTLA4  
195 (-) PD1 (+), CTLA4 (+) PD1 (-) and CTLA4 (+) PD1 (+) in LUSC (Fig.10.B-D). There  
196 was no significant correlation between IPS receiving CTLA4 (-) PD1 (-), CTLA4 (-)  
197 PD1 (+), CTLA4 (+) PD1 (-) and CTLA4 (+) PD1 (+) treatment in HNSC (Fig.10.E-  
198 H).

### 199 **Discussion**

200 Even though multiple genes (COL7A1, HJURP, KIF18B and TMEM132A) with high  
201 accuracy in the training set were calculated by LASSO regression and SVM-RFE  
202 algorithm, only COL7A1 showed high accuracy in the validation set (AUC: 0.907). In  
203 specific cancer types, COL7A1 has poor prediction performance in ESCA samples in  
204 the training set and CESC samples in the validation set. However, COL7A1 has high



205 accuracy in the diagnosis of HNSC and LUSC. Whether in the training set or the  
206 validation set, COL7A1 has the potential to be used as a diagnostic marker for HNSC  
207 and LUSC. The results showed that COL7A1 was worth exploring as a diagnostic  
208 marker for LUSC (AUC in the training set was 0.987 and AUC in the validation set was  
209 1). COL7A1 encodes collagen VII (C1), which is involved in the assembly of anchoring  
210 fibrils that fix the epidermis and dermis<sup>29</sup>. Previous studies have shown that alteration  
211 of COL7A1 can cause Recessive dystrophic epidermolysis bullosa (RDEB), and RDEB  
212 can cause skin fragility and persistent blisters, and can rapidly progress to fibrosis and  
213 even SCCs<sup>29,30</sup>. Compared with normal esophageal tissue, COL7A1 is highly expressed  
214 in ESCA. COL7A1 is significantly correlated with depth of tumor invasion and  
215 lymphatic invasion in ESCA, and the expression of COL7A1 is significantly negatively  
216 correlated with the 0-year survival rate of ESCA patients<sup>31</sup>. Although there was no  
217 significant correlation between the expression of COL7A1 and the age, gender and  
218 stage of patients including LUSC and HNSC patients, the results of prognostic analysis  
219 showed that the high expression of COL7A1 was significantly associated with shorter  
220 OS and DSS in HNSC and LUSC patients, indicating that COL7A1 also has the  
221 potential as a prognostic marker. In LUSC, the expression of COL7A1 was significantly  
222 negatively correlated with StromalScore, ImmuneScore and ESTIMATEScore, as well  
223 as the infiltration of B cell, T cell CD8, Macrophage and DC, indicating that COL7A1  
224 affects the tumor microenvironment and tumor immunity of LUSC. Further  
225 immunotherapy analysis showed that the expression of COL7A1 was significantly  
226 negatively correlated with IPS in LUSC patients treated with CTLA4 (-) PD1 (+),  
227 CTLA4 (+) PD1 (-) and CTLA4 (+) PD1 (+), indicating that COL7A1 may also have  
228 the prospect of predicting the efficacy of LUSC immunotherapy.

229 Due to the unpredictability of the results of machine learning algorithms, our study has  
230 not yet obtained appropriate clinical tissues to verify its predictive performance as a  
231 diagnostic marker. However, our study includes different databases and data sets. The  
232 results of statistical calculations on the differential expression of COL7A1 in these  
233 cancer tissues and adjacent tissues are surprising, which provides a promising marker  
234 for the early diagnosis of LUSC and HNSC.

### 235 **Conclusion**

236 COL7A1 has the potential to be a diagnostic marker for LUSC and HNSC (especially  
237 in LUSC). The high expression of COL7A1 is significantly correlated with shorter OS  
238 in HNSC patients and LUSC patients and is significantly negatively correlated with IPS

239 in LUSC patients receiving CTLA4 (-) PD1 (+), CTLA4 (+) PD1 (-) and CTLA4 (+)  
240 PD1 (+) immunotherapy.

#### 241 **Declarations**

#### 242 **Ethics approval and consent to participate**

243 Not applicable

#### 244 **Consent for publication**

245 Not applicable

#### 246 **Availability of data and materials**

247 The datasets analyzed during the current study are available in TCGA  
248 (<https://portal.gdc.cancer.gov/>), GEO database (<https://www.ncbi.nlm.nih.gov/geo/>),  
249 TISIDB database (<http://cis.hku.hk/TISIDB/>) and TCIA database (<https://tcia.at/home>).

#### 250 **Competing interests**

251 The authors declare that they have no conflicts of interest to report regarding the present  
252 study.

#### 253 **Funding**

254 Not applicable

255 **Author Contributions:** Xianglai Jiang and Chenyu Wang conceived the study,  
256 Yongxin Ma comprehensively collected relevant data, Chenyu Wang completed the  
257 work on data analysis, Xianglai Jiang completed the draft, Jiaojiao Qi and Xianglai  
258 Jiang reviewed the paper.

#### 259 **Acknowledgement**

260 This work has benefited from the aforementioned databases.

#### 261 **Reference**

- 262 1. Johnson, D. E.; Burtneß, B.; Leemans, C. R.; Lui, V. W. Y.; Bauman, J. E.; Grandis, J. R.,  
263 Head and neck squamous cell carcinoma. *Nat Rev Dis Primers* **2020**, *6* (1), 92.
- 264 2. Kallini, J. R.; Hamed, N.; Khachemoune, A., Squamous cell carcinoma of the skin:  
265 epidemiology, classification, management, and novel trends. *Int J Dermatol* **2015**, *54* (2), 130-40.
- 266 3. Small, W., Jr.; Bacon, M. A.; Bajaj, A.; Chuang, L. T.; Fisher, B. J.; Harkenrider, M. M.;  
267 Jhingran, A.; Kitchener, H. C.; Mileskin, L. R.; Viswanathan, A. N.; Gaffney, D. K., Cervical  
268 cancer: A global health crisis. *Cancer* **2017**, *123* (13), 2404-2412.
- 269 4. Codipilly, D. C.; Wang, K. K., Squamous Cell Carcinoma of the Esophagus. *Gastroenterol Clin*  
270 *North Am* **2022**, *51* (3), 457-484.
- 271 5. Wei, K. X.; Hoang, L. N., Squamous and Glandular Lesions of the Vulva and Vagina: What's  
272 New and What Remains Unanswered? *Surg Pathol Clin* **2022**, *15* (2), 389-405.
- 273 6. Lau, S. C. M.; Pan, Y.; Velcheti, V.; Wong, K. K., Squamous cell lung cancer: Current  
274 landscape and future therapeutic options. *Cancer Cell* **2022**, *40* (11), 1279-1293.
- 275 7. Waldman, A.; Schmults, C., Cutaneous Squamous Cell Carcinoma. *Hematol Oncol Clin North*  
276 *Am* **2019**, *33* (1), 1-12.

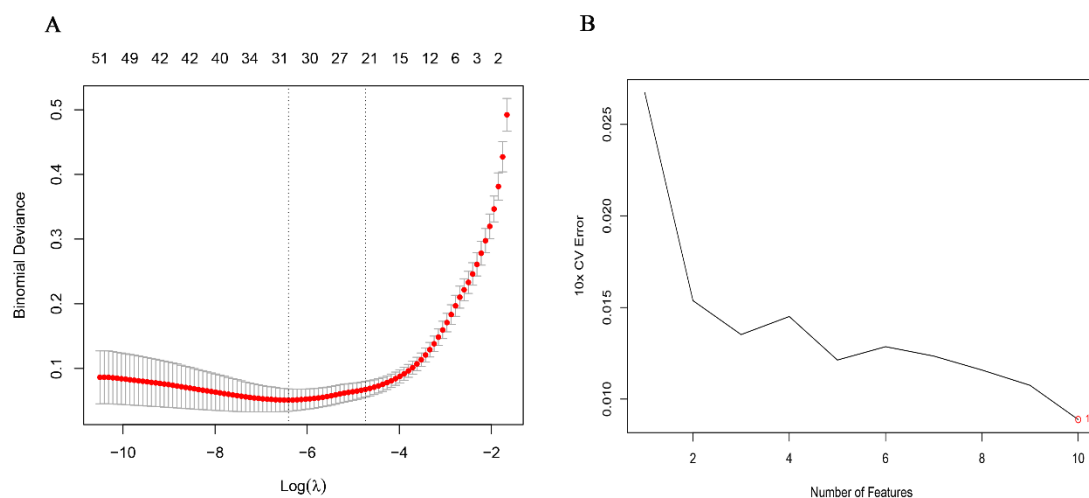


- 277 8. Derman, B. A.; Mileham, K. F.; Bonomi, P. D.; Batus, M.; Fidler, M. J., Treatment of  
278 advanced squamous cell carcinoma of the lung: a review. *Translational lung cancer research* **2015**,  
279 4 (5), 524.
- 280 9. Blobel, G. A.; Moll, R.; Franke, W. W.; Vogt-Moykopf, I., Cytokeratins in normal lung and  
281 lung carcinomas: I. Adenocarcinomas, squamous cell carcinomas and cultured cell lines. *Virchows*  
282 *Archiv B* **1984**, 45, 407-429.
- 283 10. Marur, S.; Forastiere, A. A. In *Head and neck squamous cell carcinoma: update on*  
284 *epidemiology, diagnosis, and treatment*, Mayo Clinic Proceedings, Elsevier: 2016; pp 386-396.
- 285 11. Abnet, C. C.; Arnold, M.; Wei, W.-Q., Epidemiology of esophageal squamous cell carcinoma.  
286 *Gastroenterology* **2018**, 154 (2), 360-373.
- 287 12. Ohashi, S.; Miyamoto, S. i.; Kikuchi, O.; Goto, T.; Amanuma, Y.; Muto, M., Recent  
288 advances from basic and clinical studies of esophageal squamous cell carcinoma.  
289 *Gastroenterology* **2015**, 149 (7), 1700-1715.
- 290 13. Nicol, A. F.; de Andrade, C. V.; Brusadelli, M. G.; Lodin, H. M.; Wells, S. I.; Nuovo, G. J.,  
291 The distribution of novel biomarkers in carcinoma-in-situ, microinvasive, and squamous cell  
292 carcinoma of the uterine cervix. *Annals of Diagnostic Pathology* **2019**, 38, 115-122.
- 293 14. Comprehensive genomic characterization of head and neck squamous cell carcinomas.  
294 *Nature* **2015**, 517(7536), 576-82.
- 295 15. Lin, D. C.; Hao, J. J.; Nagata, Y.; Xu, L.; Shang, L.; Meng, X.; Sato, Y.; Okuno, Y.;  
296 Varela, A. M.; Ding, L. W.; Garg, M.; Liu, L. Z.; Yang, H.; Yin, D.; Shi, Z. Z.; Jiang, Y. Y.;  
297 Gu, W. Y.; Gong, T.; Zhang, Y.; Xu, X.; Kalid, O.; Shacham, S.; Ogawa, S.; Wang, M. R.;  
298 Koeffler, H. P., Genomic and molecular characterization of esophageal squamous cell carcinoma.  
299 *Nat Genet* **2014**, 46 (5), 467-73.
- 300 16. Comprehensive genomic characterization of squamous cell lung cancers. *Nature* **2012**, 489  
301 (7417), 519-25.
- 302 17. Fukusumi, T.; Califano, J. A., The NOTCH Pathway in Head and Neck Squamous Cell  
303 Carcinoma. *J Dent Res* **2018**, 97(6), 645-653.
- 304 18. Mori, T., Involvement of the p53-p16/RB pathway control mechanism in early-stage  
305 carcinogenesis in head and neck squamous cell carcinoma. *Pathol Int* **2022**, 72(12), 577-588.
- 306 19. Shi, W.; Lee, K. E.; Wahba, G. In *Detecting disease-causing genes by LASSO-Patternsearch*  
307 *algorithm*, BMC proceedings, BioMed Central: 2007; pp 1-5.
- 308 20. Sanz, H.; Valim, C.; Vegas, E.; Oller, J. M.; Reverter, F., SVM-RFE: selection and  
309 visualization of the most relevant features through non-linear kernels. *BMC bioinformatics* **2018**,  
310 19(1), 1-18.
- 311 21. Zou, H., The adaptive lasso and its oracle properties. *Journal of the American statistical*  
312 *association* **2006**, 101 (476), 1418-1429.
- 313 22. Hou, J.; Aerts, J.; den Hamer, B.; van Ijcken, W.; den Bakker, M.; Riegman, P.; van der  
314 Leest, C.; van der Spek, P.; Foekens, J. A.; Hoogsteden, H. C.; Grosveld, F.; Philipsen, S., Gene  
315 expression-based classification of non-small cell lung carcinomas and survival prediction. *PLoS*  
316 *One* **2010**, 5 (4), e10312.
- 317 23. Scotto, L.; Narayan, G.; Nandula, S. V.; Arias-Pulido, H.; Subramaniam, S.; Schneider,  
318 A.; Kaufmann, A. M.; Wright, J. D.; Pothuri, B.; Mansukhani, M.; Murty, V. V., Identification of  
319 copy number gain and overexpressed genes on chromosome arm 20q by an integrative genomic  
320 approach in cervical cancer: potential role in progression. *Genes Chromosomes Cancer* **2008**, 47

- 321 (9), 755-65.
- 322 24. Hyland, P. L.; Zhang, H.; Yang, Q.; Yang, H. H.; Hu, N.; Lin, S. W.; Su, H.; Wang, L.;  
323 Wang, C.; Ding, T.; Fan, J. H.; Qiao, Y. L.; Sung, H.; Wheeler, W.; Giffen, C.; Burdett, L.;  
324 Wang, Z.; Lee, M. P.; Chanock, S. J.; Dawsey, S. M.; Freedman, N. D.; Abnet, C. C.;  
325 Goldstein, A. M.; Yu, K.; Taylor, P. R., Pathway, in silico and tissue-specific expression quantitative  
326 analyses of oesophageal squamous cell carcinoma genome-wide association studies data. *Int J*  
327 *Epidemiol* **2016**, *45* (1), 206-20.
- 328 25. Chen, C.; Méndez, E.; Houck, J.; Fan, W.; Lohavanichbutr, P.; Doody, D.; Yueh, B.;  
329 Futran, N. D.; Upton, M.; Farwell, D. G.; Schwartz, S. M.; Zhao, L. P., Gene expression profiling  
330 identifies genes predictive of oral squamous cell carcinoma. *Cancer Epidemiol Biomarkers Prev*  
331 **2008**, *17* (8), 2152-62.
- 332 26. Tompkins, J. D., Discovering DNA Methylation, the History and Future of the Writing on DNA.  
333 *J Hist Biol* **2022**, *55* (4), 865-887.
- 334 27. Dai, X.; Ren, T.; Zhang, Y.; Nan, N., Methylation multiplicity and its clinical values in cancer.  
335 *Expert Rev Mol Med* **2021**, *23*, e2.
- 336 28. Legrand, O.; Perrot, J.-Y.; Baudard, M.; Cordier, A.; Lautier, R. g.; Simonin, G.;  
337 Zittoun, R.; Casadevall, N.; Marie, J.-P., The immunophenotype of 177 adults with acute myeloid  
338 leukemia: proposal of a prognostic score. *Blood, The Journal of the American Society of*  
339 *Hematology* **2000**, *96* (3), 870-877.
- 340 29. Bolton, L., New options to manage epidermolysis bullosa. *Wounds* **2022**, *34* (12), 297-299.
- 341 30. Pfindner, E. G.; Lucky, A. W., Dystrophic Epidermolysis Bullosa. In *GeneReviews*(@), Adam, M.  
342 P.; Mirzaa, G. M.; Pagon, R. A.; Wallace, S. E.; Bean, L. J. H.; Gripp, K. W.; Amemiya, A., Eds.  
343 University of Washington, Seattle  
344 Copyright © 1993-2023, University of Washington, Seattle. GeneReviews is a registered trademark  
345 of the University of Washington, Seattle. All rights reserved.: Seattle (WA), 1993.
- 346 31. Kita, Y.; Mimori, K.; Tanaka, F.; Matsumoto, T.; Haraguchi, N.; Ishikawa, K.;  
347 Matsuzaki, S.; Fukuyoshi, Y.; Inoue, H.; Natsugoe, S., Clinical significance of LAMB3 and COL7A1  
348 mRNA in esophageal squamous cell carcinoma. *European Journal of Surgical Oncology (EJSO)*  
349 **2009**, *35* (1), 52-58.

350

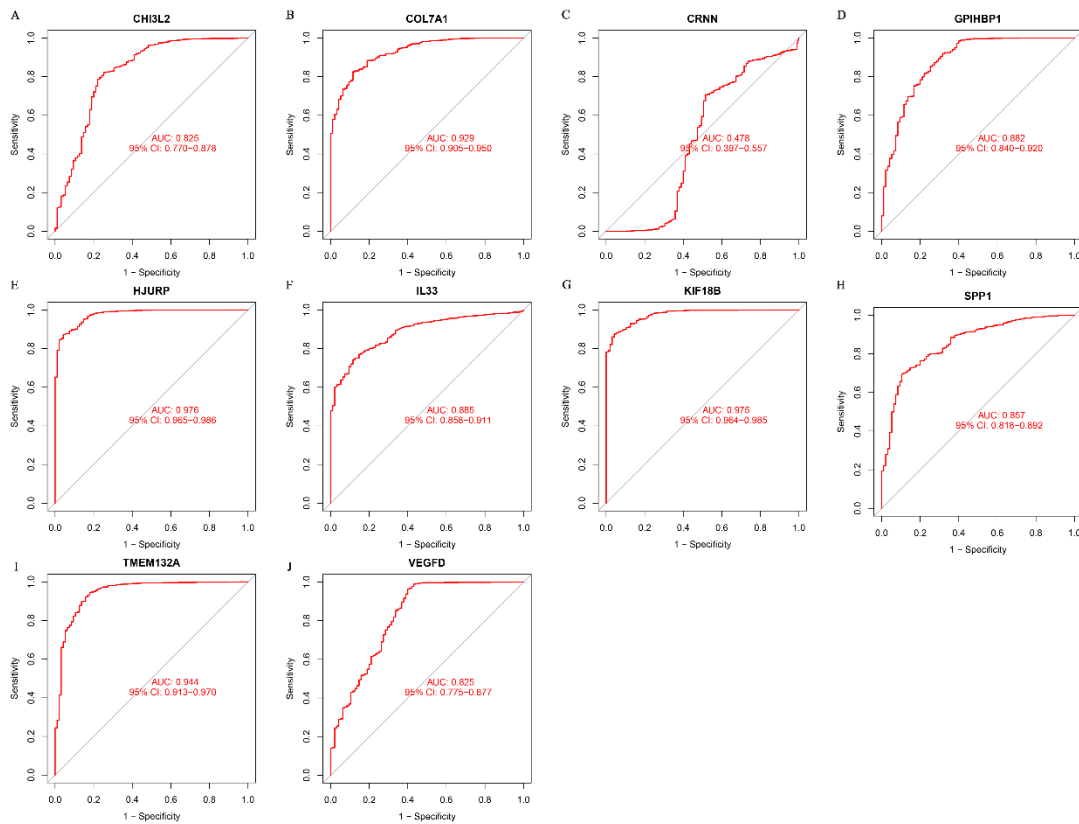
351 Figure 1: LASSO regression (A) and SVM-RFE (B) calculation results.



352

353

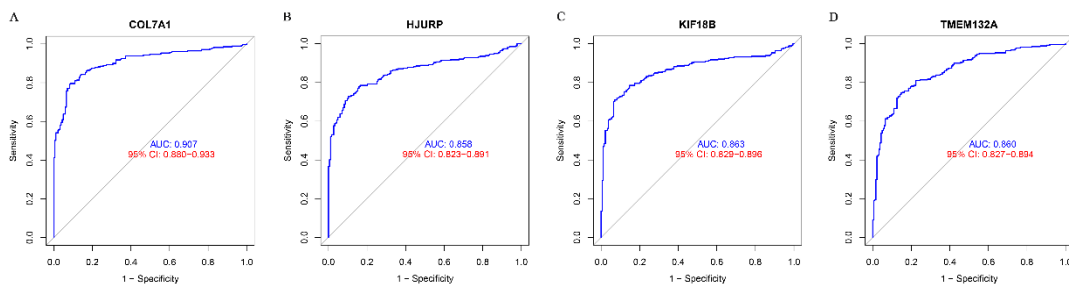
354 Figure 2: ROC curves of CHI3L2 (A), COL7A1 (B), CRNN (C), GPIHBP1 (D),  
355 HJURP (E), IL33 (F), KIF18B (G), SPP1 (H), TMEM132A (I) and VEGFD (J) in the  
356 diagnosis of SCCs in TCGA database.



357

358

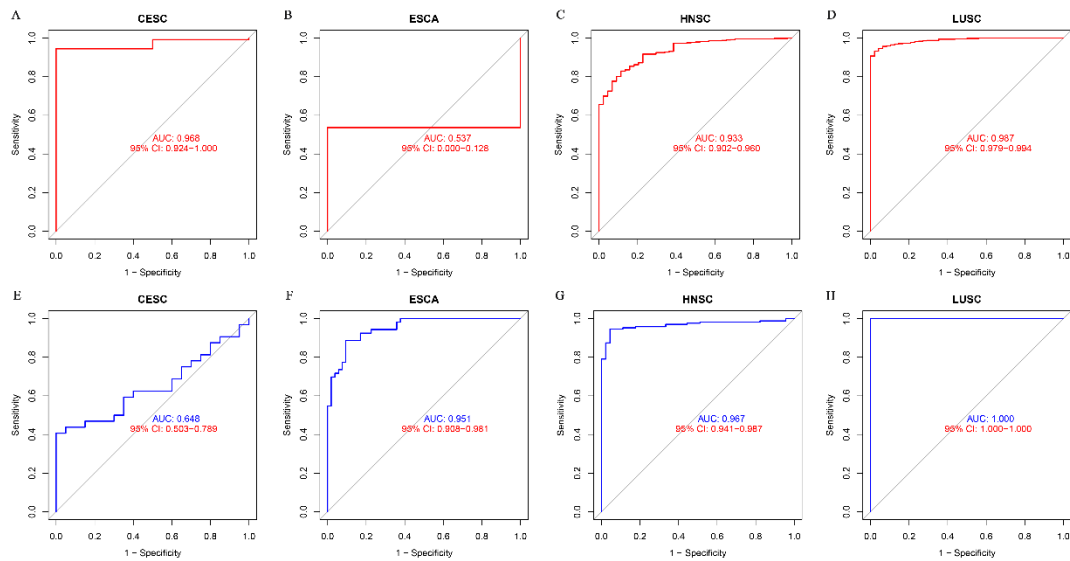
359 Figure 3: ROC curves of COL7A1 (A), HJURP (B), KIF18B (C) and TMEM132A  
360 (D) in the diagnosis of SCCs in GEO database.



361

362

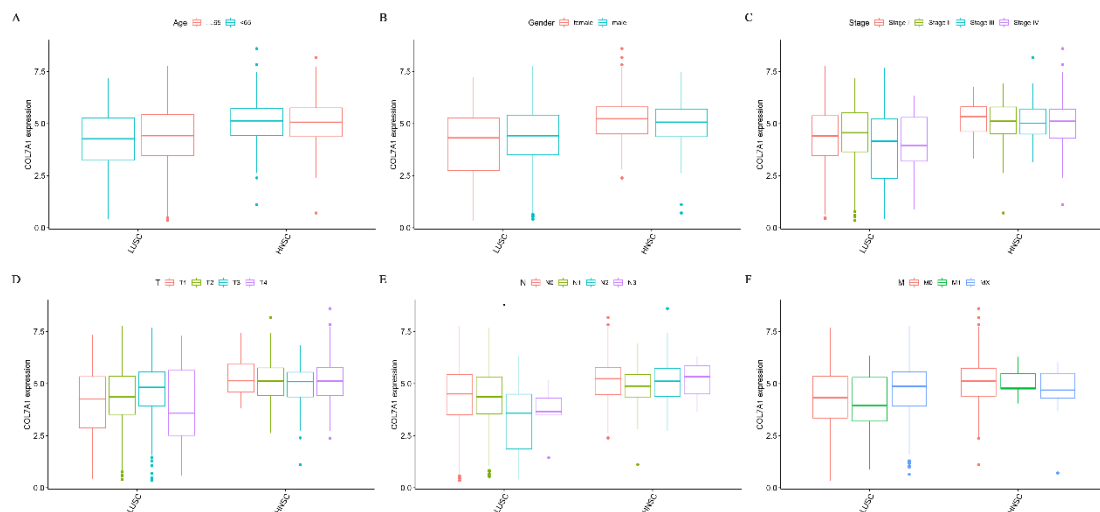
363 Figure 4: The ROC curves of COL7A1 in the diagnosis of CESC (A), ESCA (B),  
364 HNSC (C) and LUSC (D) in TCGA database; the ROC curve of COL7A1 in the  
365 diagnosis of CESC (E), ESCA (F), HNSC (G) and LUSC (H) in GEO database.



366

367

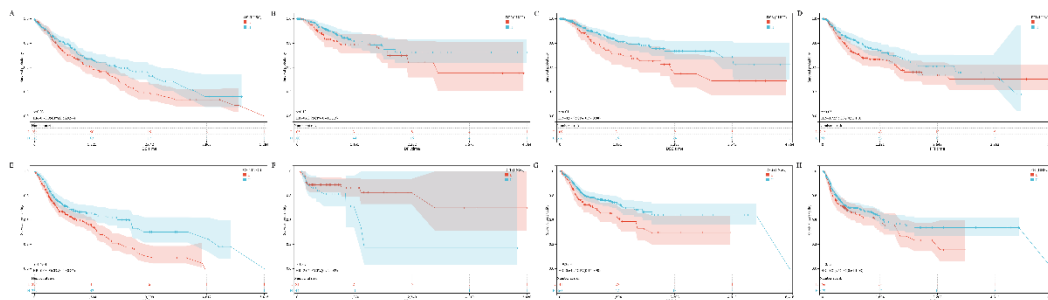
368 Figure 5: The correlation between COL7A1 expression and age (A), gender (B), stage  
369 (C), T stage (D), N stage (E) and M stage (F).



370

371

372 Figure 6: The correlation between the expression of COL7A1 and OS (A), DSS (B),  
373 DFI (C), PFI (D) in LUSC patients; the correlation between COL7A1 expression and  
374 OS (E), DSS (F), DFI (G), PFI (H) in HNSC patients.

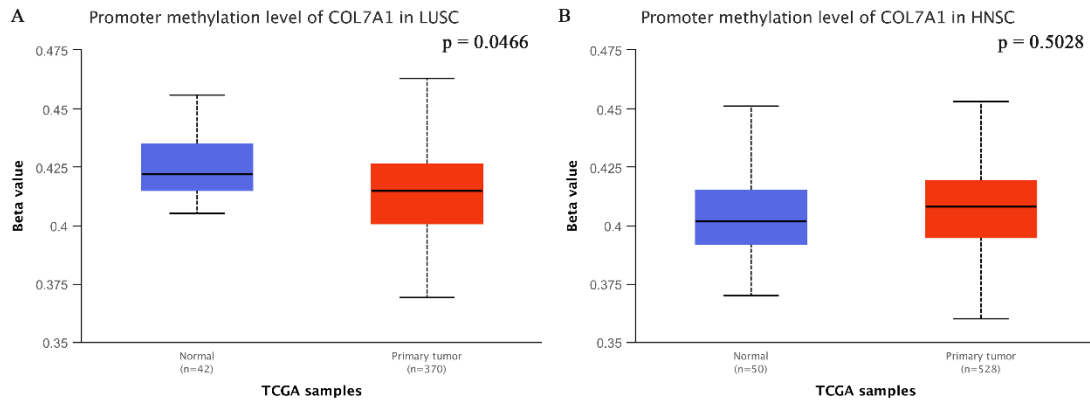


375

376

377 Figure 7: The difference of COL7A1 promoter methylation level between LUSC

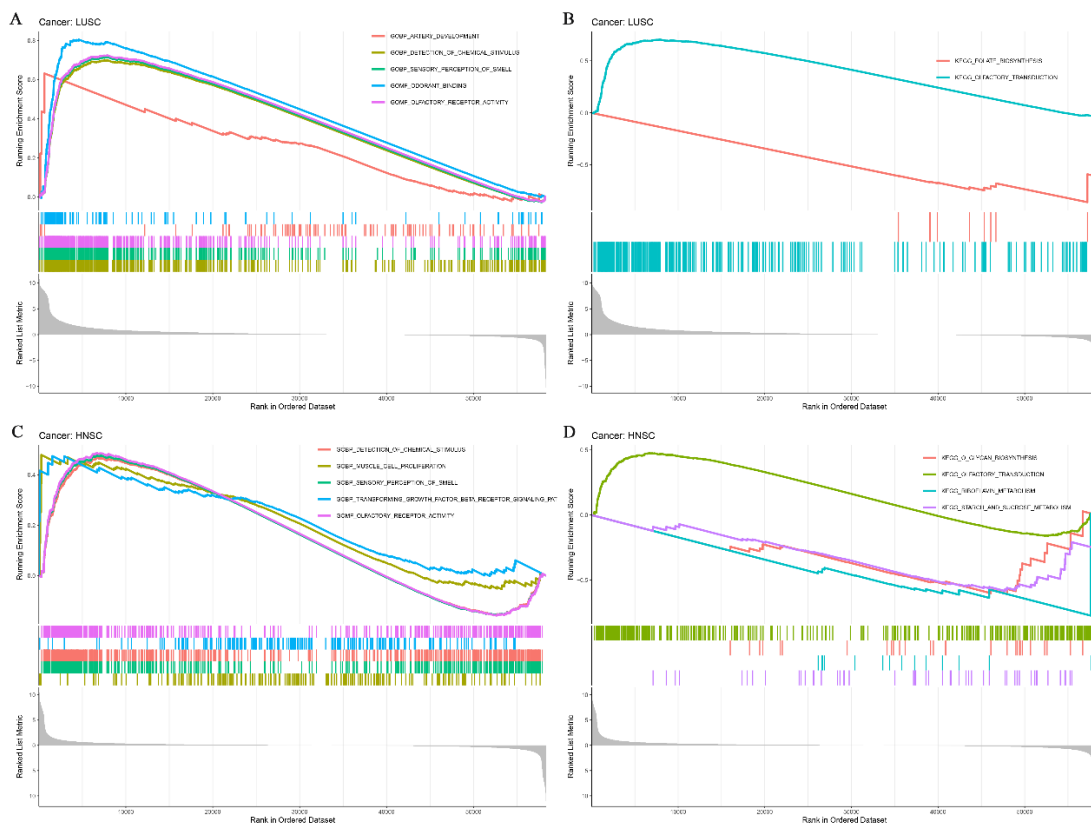
378 samples and adjacent samples (A); the difference of COL7A1 promoter methylation  
 379 level between LUSC samples and adjacent samples (B).



380

381

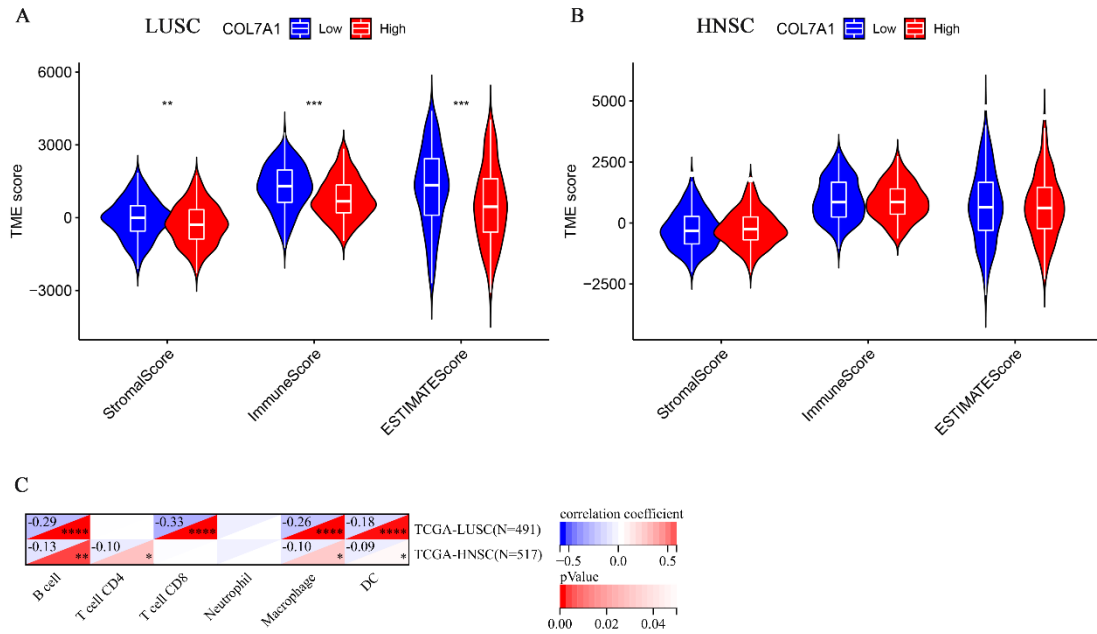
382 Figure 8: GO enrichment analysis (A) and KEGG enrichment (B) of COL7A1 in  
 383 LUSC based on GESA; GO enrichment analysis (C) and KEGG enrichment (D) of  
 384 COL7A1 in HNSC based on GESA.



385

386

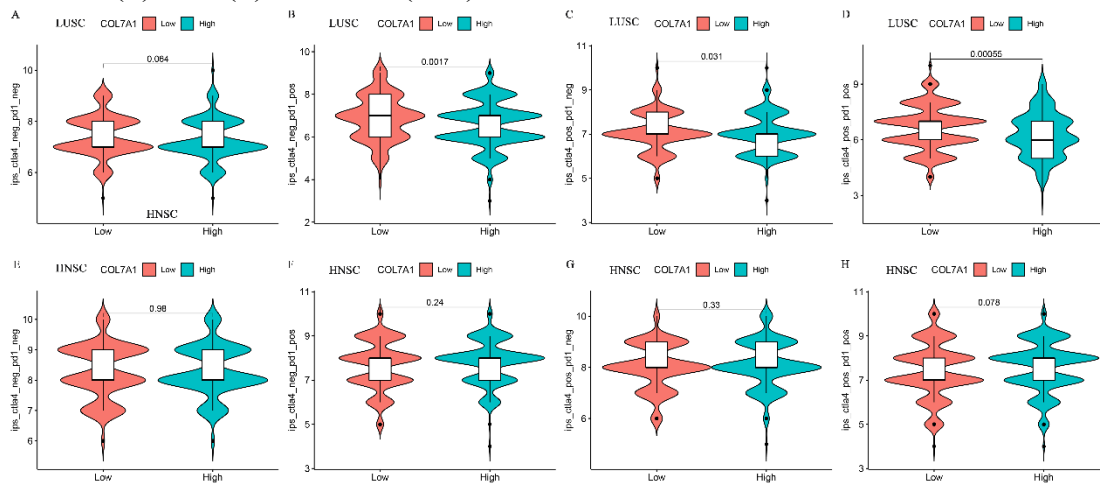
387 Figure 9: The correlation between the expression of COL7A1 and the tumor  
 388 microenvironment score in LUSC (A) ; the correlation between the expression of  
 389 COL7A1 and the tumor microenvironment score in HNSC (B) ; the correlation  
 390 between COL7A1 expression and immune cell infiltration (C).



391

392

393 Figure 10: The correlation between COL7A1 expression and IPS in LUSC patients  
 394 treated with CTLA4 (-) PD1 (-), CTLA4 (-) PD1 (+), CTLA4 (+) PD1 (-) and CTLA4  
 395 (+) PD1 (+) (A-D) ; the correlation between COL7A1 expression and IPS in HNSC  
 396 patients receiving CTLA4 (-) PD1 (-), CTLA4 (-) PD1 (+), CTLA4 (+) PD1 (-) and  
 397 CTLA4 (+) PD1 (+) treatment (E-H).



398

399