Personalized Medicine:

Redefining Cancer Treatment

Member







Steps in Data Science



Outline

O1
DATA
UNDERSTANDING

03
MODELING

05 CONCLUSION & DISCUSSION

02
DATA
PREPARATION

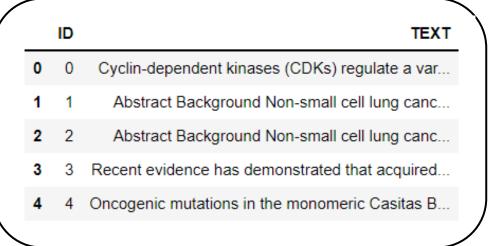
04
RESULT

01 Data Understanding

- เป็นข้อมูล genetic mutations ที่ผลมาจากผลตรวจทาง clinical ในรูปแบบ text
- ข้อมูล genetic mutation นี้ แบ่งเป็น <u>9 class*</u>
- ข้อมูลที่ใช้มีทั้งหมด 2 ไฟล์ ได้แก่
 - data_variants (ID, Gene, Variations, Class)
 - data_text (ID, Text)

	ID	Gene	Variation	Class	
0	0	FAM58A	Truncating Mutations	1	
1	1	CBL	W802*	2	
2	2	CBL	Q249E	2	
3	3	CBL	N454D	3	
4	4	CBL	L399V	4	

<mark>้ตัวอย่างข้อมูลจากไฟล์ data_variants</mark>



ตัวอย่างข้อมูลจากไฟล์ data_text

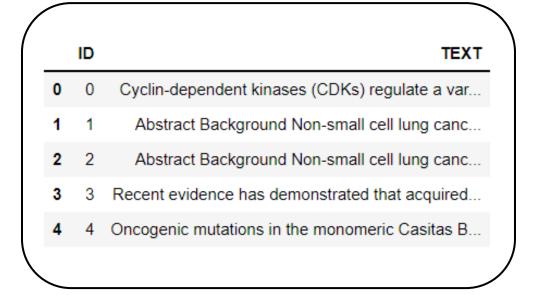
01 Data Understanding

ตัวอย่างข้อมูลจากไฟล์ data_variants

	ID	Gene	Variation	Class
0	0	FAM58A	Truncating Mutations	1
1	1	CBL	W802*	2
2	2	CBL	Q249E	2
3	3	CBL	N454D	3
4	4	CBL	L399V	4

- 1.ID = (the id of the row used to link the mutation to the clinical evidence),
- 2. **Gene** (the gene where this genetic mutation is located),
- 3. Variation (the aminoacid change for this mutations),
- 4. Class (1-9 the class this genetic mutation has been classified on)

ตัวอย่างข้อมูลจากไฟล์ data_text



- 1.ID (the id of the row used to link the clinical evidence to the genetic mutation),
- 2. **Text** (the clinical evidence used to classify the genetic mutation)

ไฟล์ data_variants

```
train_variants.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3321 entries, 0 to 3320
Data columns (total 4 columns):
              Non-Null Count Dtype
    Column
    ID 3321 non-null int64
    Gene 3321 non-null object
    Variation 3321 non-null
                              object
    Class
               3321 non-null
                              int64
dtypes: int64(2), object(2)
memory usage: 103.9+ KB
#ตรวจสอบ na ในข้อมูล
train_variants.isna().sum()
ID
Gene
Variation
Class
dtype: int64
```

ไฟล์ data_text

```
train_text.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3321 entries, 0 to 3320
Data columns (total 2 columns):
    Column Non-Null Count Dtype
             3321 non-null int64
             3316 non-null object
    TEXT
dtypes: int64(1), object(1)
memory usage: 52.0+ KB
#ตรวจสอบ na ในข้อมูล
train_text.isna().sum()
ID
        <sub>5</sub> มี na ในข้อมูลอยู่ 5
TEXT
dtype: int64
```

• จัดการกับ na ในไฟล์ data_text

```
#จัดการกับ na ในข้อมูล โดยเติม unknown เข้าไป
train_text['TEXT'].fillna('unknown',inplace=True)

#ตรวจสอบ na ในข้อมูล
train_text.isna().sum()

ID 0
TEXT 0 ไม่มี na ในข้อมูลแล้ว
dtype: int64
```

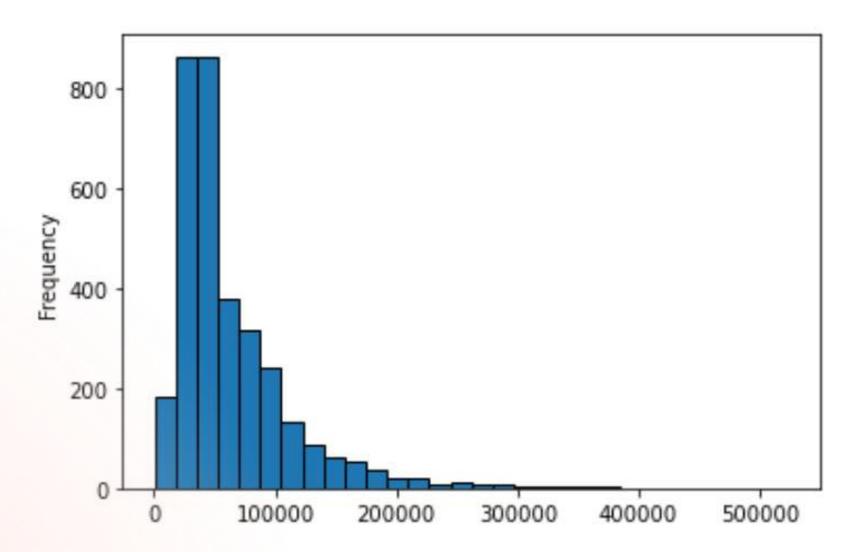
data_variants + data_text = data_merge

```
[9] # merge ไฟล์ทั้งสอง
train_merge = pd.merge(train_variants, train_text, on='ID')
train_merge.head()
```

	ID	Gene	Variation	Class	TEXT
0	0	FAM58A	Truncating Mutations	1	Cyclin-dependent kinases (CDKs) regulate a var
1	1	CBL	W802*	2	Abstract Background Non-small cell lung canc
2	2	CBL	Q249E	2	Abstract Background Non-small cell lung canc
3	3	CBL	N454D	3	Recent evidence has demonstrated that acquired
4	4	CBL	L399V	4	Oncogenic mutations in the monomeric Casitas B

• ความยาวของ Text ในข้อมูล

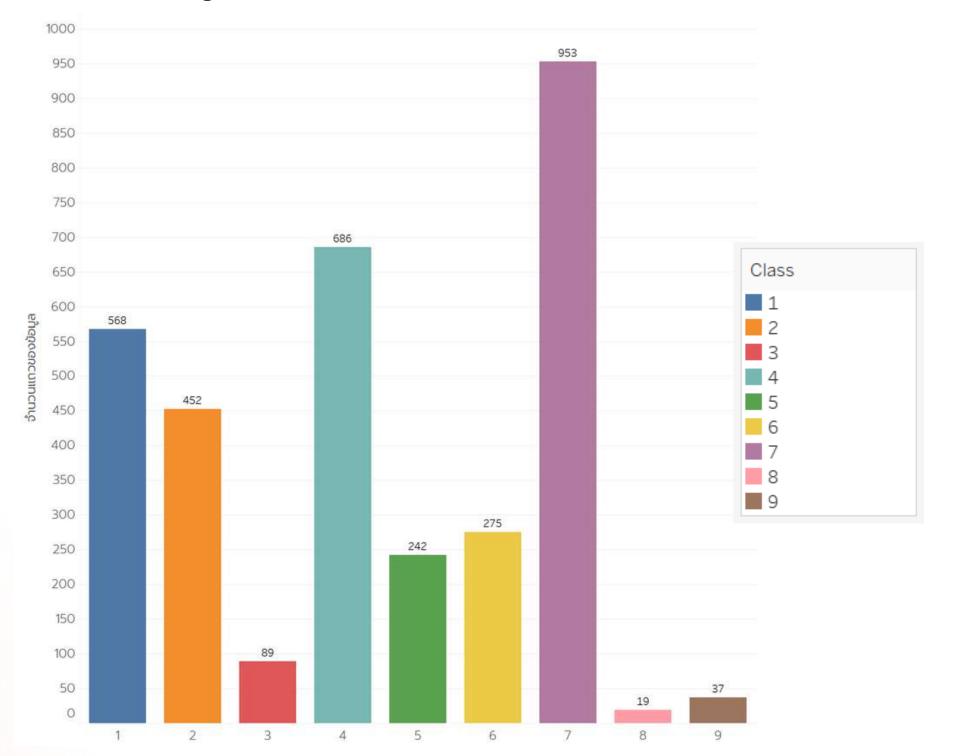
```
1 # เพิ่ม col length
2 train_merge['length'] = train_merge['TEXT'].str.len()
```



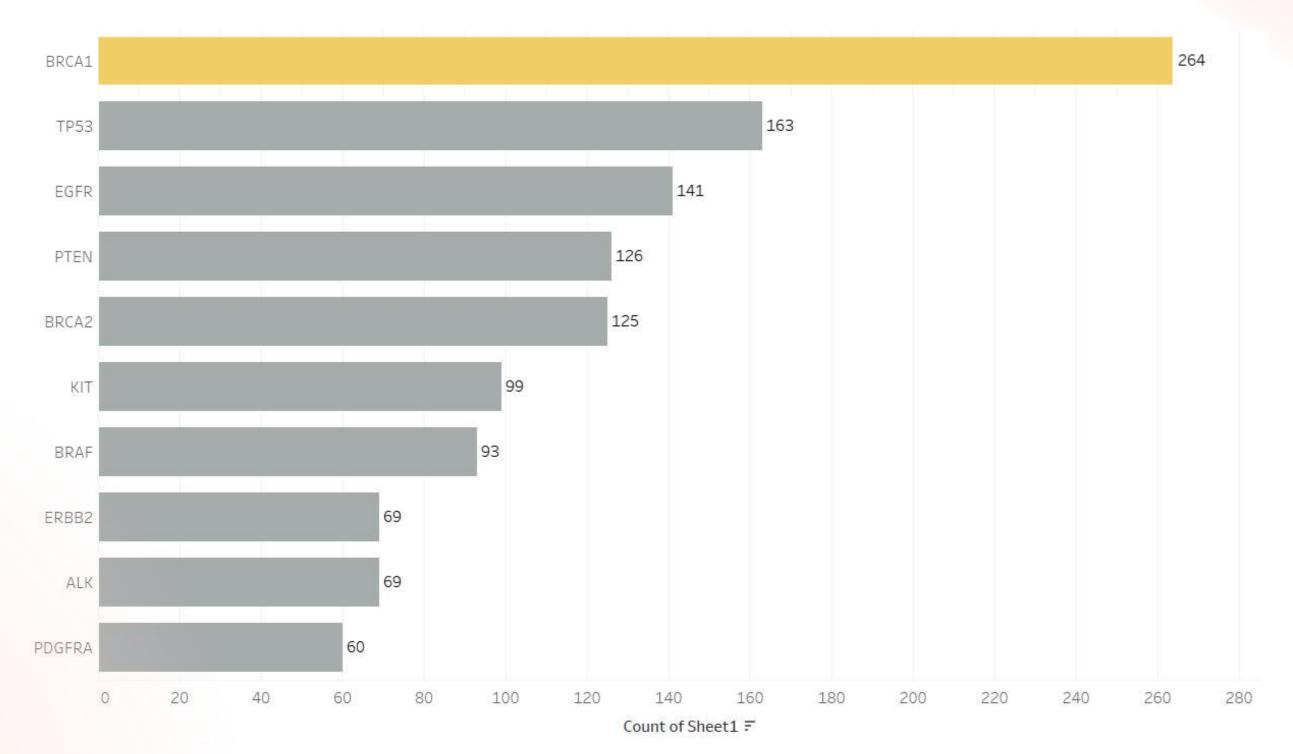
ตัวอย่างข้อมูลจาก column TEXT เพียง 1 row โดยมีความยาว (str.len()) อยู่ที่ 39672

'Cyclin-dependent kinases (CDKs) regulate a variety of fundamental cellular processes. CDK10 stands out as one of the last orphan CDKs for which no activation of the MAPK pathway, which confers tamoxifen resistance to breast cancer cells. The precise mechanisms by which CDK10 modulates ETS2 (v-ets enythroblastosis virus E26 oncogene homolog 2)-driven activation of the MAPK pathway, which confers tamoxifen resistance to breast cancer cells. The precise mechanisms by which CDK10 modulates ETS2 (v-ets enythroblastosis virus EA6 oncogene homolog 2)-driven activation of the MAPK pathway, which confers tamoxifen resistance to breast cancer cells. The precise mechanisms by which CDK10 modulates ETS2 activity, and more generally the functions of the MAPK pathway, which confers tamoxifen resistance to breast cancer cells. The precise mechanisms by which CDK10 modulates ETS2 activity, and more generally the functions of the MAPK pathway. 0, remain elusive. Here we demonstrate that CDK10 is a cyclin-dependent kinase by identifying cyclin M as an activating cyclin M as an activating cyclin M, an orphan cyclin, is the product of FAM58A, whose mutations cause STAR syndrome, a human developmental anomable to interact with CDK10. Cyclin M silencing phenocopies CDK10 in M remains and in conferring in increasing c-Raf and in conferring wifen resistance to breast cancer cells. CDK10/cyclin M phosphorylates ETS2 in vitro, and in cells it positively controls ETS2 degradation by the proteasome. ETS2 protein levels are increased in cells derived from a STAR syndrome. Cyclin M levels. Altogether, our results reveal an additional regulatory mechanism for ETS2, which plays key roles in cancer and development. They also shed light on the molecular mechanisms underlying STAR syndrome. Cyclin dependent kinases (CDK3) play a pivotal role in the control of a number of fundamental cellular processes (1). The human genome contains 21 genes encoding proteins that can be considered as members of the CDK family owing to their sequence similarity with bona fide CDKs, those known to be activated by cyclins (2). Although discovered almost 20 y ago (3, 4), CDK10 remains one of the two CDKs without an identified cyclin partner. This knowledge gap has largely impeded the exploration of its biological functions. CDK10 can act as a positive cell cycle regulator in some cells (5, 6) or as a tumor suppressor in others (7, 8). CDK10 interacts with the ETS2 (v-ets erythroblastosis virus E26 oncogene homolog 2) transcription factor and inhibits its transcriptional activity through an unknown mechanism (9). CDK10 knockdown derepresses ETS2, which increases the expression of the c-Raf protein kinase, activates the MAPK pathway, and induces resistance of MCF7 cells to tamoxifen (6). Here, we deorphanize CDK10 by identifying cyclin M, the product of FAM58A, as a binding partner. Mutations in this gene that predict absence or truncation of cyclin M are associated with STAR syndrome remain unknown. We show that a recombinant CDK10/cyclin M heterodimer is an active protein kinase that phosphorylates ETS2 in vitro. Cyclin M silencing phenocopies CDK10 silencing in increasing c-Raf and phosphor-ERK expression levels and in inducing tamoxifen resistance in estrogen receptor (ER)+ breast cancer cells. We show that CDK10/cyclin M posprosed in the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived from a STAR patient, and we demonstrate that it is attributable to the decreased cyclin M expression level in cells derived fro syndrome patients (10). None of these shorter isoforms produced interaction phenotypes (Fig. 1.A and C and Fig. S1A). Fig. 1.In a new window Download PPTFig. 1.CDK10 and cyclin M form an interaction complex. (A) Schematic representation of the different protein isoforms analyzed by YZH assays. Amino acid numbers are indicated. Black boxes indicate internal deletions. The red box indicates a differing amino acid sequence compared with CDK10 P1. (B) YZH assay between a set of CDK proteins expressed as baits (in fusion to the LexA DNA binding domain) and CDK interacting proteins expressed as preys (in fusion to the B42 transcriptional activator). pEG202 and pJG4-5 are the empty bait and prey plasmids expressing LexA and B42, respectively. lacZ was used as a reporter gene, and blue yeast are indicative of a Y2H interaction phenotype. (C) Y2H assay between the different protein interaction phenotype. (C) Y2H assay between the different protein is protein and complex. (A) Schematic representation of the different protein is protein and cyclin M form an interaction phenotype. (E) Y2H assay between a set of CDK protein analyzed by Y2H assay, Amino acid numbers are indicated. Black boxes indicate internal deletions. The red box indicates a differing amino acid sequence compared with CDK10 P1. (B) Y2H assay between a set of CDK protein sexpension of the different protein is protein sexpension. (B) Y2H assay between a set of CDK protein sexpension of the different protein is protein sexpension. (B) Y2H assay between a set of CDK protein sexpension of the different protein is protein sexpension. (B) Y2H assay between a set of CDK protein sexpension of the different protein is protein sexpension. (B) Y2H assay between a set of CDK protein sexpension of the different protein se 6His expression levels in transfected HEK293 cells. (E) Western blot analysis of Myc-CDK10 (wt or kd) immunoprecipitates obtained using the anti-Myc antibody, "Inputs" corresponds to 30 µg MCF7 total cell lysates. The lower band of the doublet observed on the upper panel comigrates with the exogenously expressed untagged CDK10 and thus corresponds to a nonspecific signal, as demonstrated by it insensitivity to either overexpression of CDK10 (Fig. S2B). Another experiment with a longer gel migration is shown in Fig. S1D.Next we examined the ability of CDK10 and cyclin M to interact when expressed in human cells (Fig. 1 D and E). We tested wild-type CDK10 (wt or kd) in a human embryonic kidney cell line (HEK293). The expression with Myc-CDK10 (wt or kd) and, to a lesser extent, that of Myc-CDK10 (wt or kd) was increased upon coexpression with Myc-CDK10 proteins and detected the presence of cyclin M in the CDK10 (wt) and (kd) immunoprecipitates only when these proteins were coexpressed pair-wise (Fig. 15). To detect the interaction between the CDK10, To detect the interaction between endogenous proteins, we performed immunoprecipitates only when these observations by detecting the presence of Myc-CDK10 in cyclin M. (Fig. 51B). These experiments confirmed these observations by detecting the presence of Myc-CDK10 in cyclin M. (Fig. 51B). These experiments confirmed the lack of robust interaction between endogenous proteins, we performed immunoprecipitates only when these observations by detecting the presence of Myc-CDK10 in cyclin M. (Fig. 51B). These experiments confirmed the lack of robust interaction between endogenous proteins by detecting the presence of Myc-CDK10 in cyclin M. (Fig. 51B). These experiments confirmed the lack of robust interaction between endogenous proteins by detecting the presence of Myc-CDK10 in cyclin M. (Fig. 51B). These experiments confirmed the lack of robust interaction between the CDK10. In cyclin M. In immunoprecipitates only when the CDK10 in cyclin M. In immunoprecipitates only when the CDK10 in cyclin M. In immunoprecipitates only when the CDK10 in cyclin M. In immunoprecipitates only when the CDK10 in cyclin M. In immunoprecipitates on the CDK10 in cyclin M. In immunopreci (Fig. 1C), is a phosphorylation substrate of CDK10/cyclin M. We detected strong phosphorylation of ETS2 by the GST-CDK10(kd)/Strepli-cyclin M purified on a glutathione Septances matrix to capture GST-CDK10(kd)/Strepli-cyclin M. Unified protein kinase. (A) In vitro protein kinase assay on histone H1. Lysates from insect cells expressing different proteins were purified on a glutathione Septances matrix to capture GST-CDK10(kd)/Strepli-cyclin M. purified protein kinase assay on histone H1. Lysates from insect cells expressing different proteins were purified on a glutathione Septances matrix to capture GST-CDK10(kd)/Strepli-cyclin M. purified protein kinase assay on histone H1. Lysates from insect cells expressing different proteins were purified on a glutathione Septance matrix to capture GST-CDK10(kd)/Strepli-cyclin M. purified protein kinase assay on histone H1. Lysates from insect cells expressing different proteins were purified on a glutathione Septance matrix to capture GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no phosphorylation was detected using GST-CDK10(kd)/Strepli-cyclin M. beterodimer, whereas no investigated whether cyclin M is also involved in this regulatory pathway. To aim at a highly specific silencing, we used siRNA pools (mix of four different siRNAs) at low final concentration (10 nM). Both CDK10 protein level (Fig. 3A and C and Fig. 52B). These results, and those shown in Fig. 10, suggest that cyclin M silencing we used siRNA pools (mix of four different siRNAs) at low final concentration (10 nM). Both CDK10 and cyclin M silencing, we used siRNA pools (mix of four different siRNAs) at low final concentration (10 nM). Both CDK10 and cyclin M silencing, we used siRNA pools (mix of four different siRNAs) at low final concentration (10 nM). Both CDK10 in a results, and those shown in Fig. 10, suggest that cyclin M silencing, we used siRNA pools (mix of four different siRNAs) at low final concentration (10 nM). Both CDK10 in a results, and those shown in Fig. 10, suggest that cyclin M silencing, we used siRNA pools (mix of four different siRNAs) at low final concentration (10 nM). Both CDK10 in a results, and those shown in Fig. 10, suggest that cyclin M silencing, we used siRNA pools (mix of four different siRNAs) at low final concentration (10 nM). Both CDK10 in a results, and those shown in Fig. 10, suggest that cyclin M silencing, we used siRNA pools (mix of four different siRNAs) at low final concentration (10 nM). Both CDK10 in mix of four different siRNAs at low final concentration (10 nM). Both CDK10 in mix of four different siRNAs at low final concentration (10 nM). Both CDK10 in mix of four different siRNAs at low final concentration (10 nM). Both CDK10 in mix of four different siRNAs at low final concentration (10 nM). Both CDK10 in mix of four different siRNAs at low final concentration (10 nM). Both CDK10 in mix of four different siRNAs at low final concentration (10 nM). Both CDK10 in mix of four different siRNAs at low final concentration (10 nM). Both CDK10 in mix of four different siRNAs at low final concentration (10 nM). Both CDK10 in mix of four different siRNAs at low final conc 0.01; ***P ≤ 0.001.We then wished to explore the mechanism by which CDK10/cyclin M controls ETS2 is a short-lived protein degraded by the proteasome (13). A straightforward hypothesis is that CDK10/cyclin M positively controls ETS2 degradation. We thus examined the impact of CDK10 and that of cyclin M caused an increase in the expression levels of an expersion levels of an expersion levels of the endogenous expersion levels. The silencing of CDK10 and that of cyclin M caused an increase in the expression levels of an expersion levels of the Flag-tagged ETS2 protein (Fig. S4B), as well as of the endogenous expersion levels. The silencing of CDK10 and that of cyclin M silencing of CDK10 and that of cyclin M silencing of CDK10 and that of cyclin M controls ETS2 degradation. We thus examined the expression levels of the Flag-tagged ETS2 protein (Fig. S4B), which controls ETS2 degradation. We thus examined the expression levels of the Flag-tagged ETS2 protein (Fig. S4B). We then examined the expression levels of the Flag-tagged ETS2 protein (Fig. S4B). We then examined the expression levels of the Flag-tagged ETS2 protein (Fig. S4B). We then examined the expression levels of the Flag-tagged ETS2 protein (Fig. S4B). We then examined the expression levels of the Flag-tagged ETS2 protein (Fig. S4B). However, its expression levels was dramatically decreased when expressed alone or in combination with Myc-CDK10 or -CDK10 o of with CNN of any Cyclin M (reg. 4.9). These occupations suggest that enloyed contracts of contract with the Map; Test and the year of the Mistract with the Map; Test and the Winthow of Winth Mistracting. A 18, the Mistracting of the Winthow of Winth Mistracting. A 18, the Mistracting of the Winthow of Winth Mistracting. A 18, the Winthow of Winth Mistracting. A 18, the Winthow of control lymphoblastoid cell line. In line with our preceding observations, we detected an increased expression level of ETS2 protein in the STAR cell line exmonstrate that the increase in ETS2 protein in the exmost of the control (Fig. 5B). To demonstrate that the increase in ETS2 protein expression level of ETS2 protein expression level of ETS2 protein many and the expression level of ETS2 mRNA similar to that of the control cell line (Fig. 5B). To demonstrate that the increase in ETS2 protein expression level of ETS2 protein expression is indeed a result of the decreased cyclin M expression observed in the STAR patient-derived lymphoblastoid cell line, we expressed cyclin M expression in STAR patient-derived lymphoblastoid cell line, derived from a healthy individual. A quantification is shown in Fig. S10A. (B) Quantitative RT-PCR analysis of cyclin M and ETS2 protein levels in the STAR patient-derived lymphoblastoid cell line, derived from a healthy individual. A quantification is shown in Fig. S10A. (B) Quantitative RT-PCR analysis of cyclin M and ETS2 protein levels in the same cells. ***P\$ = 0.001. (C) Western blot revealing endogenously expressed cyclin M levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels is shown in Fig. S10B. A quantification of ETS2 protein levels in the STAR patient-derived lymphoblastoid cell line, derived lymphoblastoid cell line, derived lymphoblastoid cell line and in a control lymphoblastoid cell line, derived lymphoblastoid cell line, deri this work, we unveil the interaction between CDX10, the last orphan CDK discovered in the pregenomic era (2), and cyclin M, the only cyclin associated with a human genetic disease so far, and whose functions remain unknown (10). The closest paralogs of the CDK10 within the CDK family are the CDK10 in the interact orbustly with cyclin M (Fig. 511). The fact that none of the shorter CDK10 in the interaction between CDK10 must be required by different observations. Both proteins, which interact robustly with cyclin M suggests that cyclin M (Fig. 512), and from the much reduced endogenous CDK10 gene (16, 17) plays an important role in regulating CDK10 in mediated degradation (18). Our observations suggest that cyclin M silencing (Fig. 3A and Fig. 52B). CDK10 is subject to ubiquitin-mediated degradation (18). Our observations suggest that cyclin M stability is enhanced upon binding to CDK10, independently from its kinase activity, as seen for cyclin CDK10, independent CDK10 proteins, which interact robustly with in the CDK1 and cyclin M silencing (Fig. 3A and Fig. 52B). CDK10 is subject to ubiquitin-mediated degradation (18). Our observations suggest that cyclin M stability is enhanced upon binding to CDK10, independently from its kinase activity, as seen for cyclin C and CDK8 (19). We uncover a cyclin M stability is enhanced upon binding to CDK10, independently from its kinase activity, as seen for cyclin C and CDK8 (19). We uncover a cyclin M stability is enhanced upon binding to CDK10, independently from its kinase activity, as seen for cyclin C and CDK8 (19). We uncover a cyclin M stability is enhanced upon binding to CDK10, independently from its kinase activity, as seen for cyclin C and CDK8 (19). We uncover a cyclin M stability is enhanced upon binding to CDK10, independently from its kinase activity in vitro, thus demonstrating that this protein with a finite cyclin M stability is enhanced upon binding to CDK10, independently from its kinase activity in vitro, thus demonstrating that this protein with a fini dependent kinase. Our Y2H assays reveal that truncated cyclin M proteins, females affected by the STAR syndrome associated FAM58A mutations do not products of two STAR syndrome associated FAM58A mutations do not products of two STAR syndrome associated FAM58A mutations do not product so five first to truncated cyclin M proteins, females affected by the STAR syndrome associated FAM58A mutations do not product so five five to truncated cyclin M proteins, females affected by the STAR syndrome associated FAM58A mutations do not product so five five to truncated cyclin M proteins, females affected by the STAR syndrome associated FAM58A mutations do not product so five five five form the decreased cyclin M proteins, females affected by the STAR syndrome associated FAM58A mutations do not product so five five five form the decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activity of ETS2 levels to partially decreased cyclin M winds activities activities activities activities activities activities activities which CDK10 may belong; Fig. S11) have been shown to phosphorylate a variety of motifs in a non-proline-directed fashion, especially in the context of molecular docking mediated phosphorylate a variety of motifs in a non-proline-directed fashion, especially in the control of ETS2 degradation involves a number of players, including APC-Cdh1 (13) and the Pointed domain of ETS2 (6, 9) (Fig. 1C) would allow docking mediated phosphorylate a variety of motifs in a non-proline-directed fashion, especially in the control of ETS2 degradation involves a number of players, including APC-Cdh1 (13) and the Pointed domain of ETS2 (9, 9) (Fig. 1C) would allow docking mediated phosphorylation of etypical sites. The control of ETS2 degradation involves a number of players, including APC-Cdh1 (13) and the Pointed domain of ETS2 (9, 9) (Fig. 1C) would allow docking mediated phosphorylation of etypical sites. The control of ETS2 degradation involves a number of players, including APC-Cdh1 (13) and the Pointed domain of ETS2 (9, 9) (Fig. 1C) would allow docking mediated phosphorylation of etypical sites. The control of ETS2 degradation involves a number of players, including APC-Cdh1 (13) and the Pointed domain of ETS2 (9, 9) (Fig. 1C) would allow docking mediated phosphorylation of etypical sites. The control of ETS2 degradation involves a number of players, including APC-Cdh1 (13) and the Pointed domain of ETS2 (9, 9) (Fig. 1C) would allow docking mediated phosphorylation of etypical sites. The control of ETS2 degradation involves a number of players, including APC-Cdh1 (13) and the cultivation of etypical sites. The control of ETS2 degradation involves a number of players, including APC-Cdh1 (13) and the cultivation of etypical sites. The control of ETS2 degradation involves a number of players, including APC-Cdh1 (13) and the cultivation of etypical sites. The control of ETS2 degradation involves a number of etypical sites. The control of ETS2 degradation involves a number of etypical sites. The control of etypical sites and et associated with reduced CDK10 expression levels, and they suggest that, like CDK10 (6), cyclin M could also be a predictive clinical marker of hormone therapy response of ERa-positive breast cancer patients. CDK10 expression of Ets2 present severe cranial abnormalities (24), and those observed in STAR patients could thus be caused at least in part by increased ETS2 protein levels. Another expected consequence of enhanced ETS2 expression levels would be a decreased risk to develop certain types of cancers and an increased risk to develop others. Studies on various mouse models (including models of Down syndrome, in which three copies of ETS2 exist) have revealed that ETS2 expression expression expression expression expression expression expression of Ets2 present severe cranial abnormalities (24), and those observed in STAR patients could thus be caused at least in part by increased ETS2 protein levels. Another expected consequence of enhanced ETS2 expression levels would be a decreased risk to develop certain types of cancers and an increased risk to develop others. Studies on various mouse models (including models of Down syndrome. Ets2 transgenic mice expositive breast cancer (25). Intringuingly, one of the very few STAR patients identified so far has been diagnosed with a nephroblastoma (26). Finally, our findings will facilitate the general exploration of the biological functions of CDK10 and, in particular, its role in the control of cell division. Previous studies have suggested either a positive role in cell cycle control (5, 6) or a tumor-suppressive activity in some cancers (7, 8). The severe growth retardation exhibited by STAR patients strongly suggests that CDK10/cyclin M plays an important role in the control of cell proliferation. Previous studies have suggested either a positive role in cell cycle control (5, 6) or a tumor-suppressive activity in some cancers (7, 8). The severe growth retardation exhibited by STAR patients strongly suggests that CDK10/cyclin M plays an important role in the contr cyclin McDAs, plasmid constructions, tamoxifen response analysis, quantitative RT-PCR, mass spectrometry experiments, and antibody production are detailed in SI Materials and Methods, Yeast Two-Hybrid Interaction Assays. We performed yeast interaction Assays. We performed yeast interaction mating assays as previously described (27). Mammalian Cell Cultures and Transfection assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection mating assays as previously described (27). Mammalian Cell Cultures and Transfection and McF7 cells in DMEM supplemented with 15% (vol/vol) FBS. (Invitrogen), and we grew lymphoblastoid cells in RPMI 15% (invitrogen) for siRNAs, and McF7 cells in DMEM supplemented with 15% (vol/vol) FBS. (Invitrogen) for siRNAs, and McF7 cells in DMEM supplemented with 15% (vol/vol) FBS. (Invitrogen) for siRNAs, and McF7 cells in DMEM supplemented with 15% (vol/vol) FBS. (Invitrogen) for siRNAs, and McF7 cells in DMEM supplemented with 15% (vol/vol) FBS. (Invitrogen) for siRNAs, an immunoprecipitation experiments on 500 µg of total proteins, in lysis buffer. We precleared the beads by scattaring wheel. We collected the beads by scat at "C and washed the beads by scattaring wheel. We collected the beads by scattaring the beads by scattaring the beads by scattaring wheel. We collected the beads by scattaring the beads by scattaring the beads by scattaring wheel. We collected the beads by scattaring the beads by scatta revealed the blots by enhanced chemiliminescence (SuperSignal West Femto, Through Circle of the cells with StrepIl-CycM—producing viruses, or coinfected the cells with StrepIl-CycM—producing viruses, and we collected the cells with StrepIl-CycM. We generated recombinant barmious in DH10Bac Escherichia coil and baculoviruses in Sf9 cells with StrepIl-CycM. We generated recombinant barmious in DH10Bac Escherichia coil and baculoviruses in Sf9 cells with StrepIl-CycM—producing viruses, or coinfected the cells with StrepIl-CycM. We generated recombinant barmious in DH10Bac Escherichia coil and baculoviruses in Sf9 cells with StrepIl-CycM. We generated recombinant barmious in DH10Bac Escherichia coil and baculoviruses in Sf9 cells with StrepIl-CycM—producing viruses, or coinfected the cells with StrepIl-CycM—producing viruses, or coinfected the cells with StrepIl-CycM—producing viruses, or coinfected the cells with StrepIl-CycM. We generated recombinant barmious in DH10Bac Escherichia coil and baculoviruses in Sf9 cells with StrepIl-CycM—producing viruses, or coinfected the cells with StrepIl-CycM—producing viruses, o the beads in 100 µL kinase buffer A containing 10% (vol/vol) glycerol for storage.6His-ETS2. We transformed Origami? DE3 (Novagen) with the 6His-ETS2 we spun 50 mL cells and resuspended the pellet in 2 mL lysis buffer (PBS, 300 mM NaCl, 10 mM Imidazole, 1 mM DTT, and 0.1% Nonidet P-40). Containing a protease inhibitor mixture without EDTA (Roche). We lysed the pellet in 2 mL lysis buffer (PBS, 300 mM NaCl, 10 mM DTT, and 0.1% Nonidet P-40). Containing a protease inhibitor mixture without EDTA (Roche). We lysed the pellet in 2 mL lysis buffer (PBS, 300 mM NaCl, 10 mM DTT, and 0.1% Nonidet P-40). Containing a protease inhibitor mixture without EDTA (Roche). We lysed the pellet in 2 mL lysis buffer (PBS, 250 mM midazole, 1 mM DTT, and 0.1% Nonidet P-40). Containing a protease inhibitor mixture without EDTA (Roche). We lysed the pellet in 2 mL lysis buffer (PBS, 250 mM not pellet in 2 mL lysis buffer (PBS, 250 m Commassie (R-250, Bio-Rad), dried then, and detected the incorporated radioactivity by autoractive lands unless provided previously in a mother-design with appearance of this syndrome. (a-f) Facial and molecular characterization of STAR syndrome. (a-f) Facial appearance 50.7 kb, with one probe positioned within FAM58A. The deletion found in case 3, which removes exon 5 and some 3' sequence. The pink horizontal bars above the boxes the horizontal line below exon 5 indicates the deletion found in case 3, which removes exon 5 and some 3' sequence. The pink horizontal bars above the boxes 1 and 2, whereas the horizontal line below exon 5 indicates the deletion found in case 3, which removes exon 5 and some 3' sequence. The pink horizontal arrow indicates the deletion found in case 3, which removes exon 5 and some 3' sequence. The pink horizontal arrow indicates the deletion found in case 3, which removes exon 5 and some 3' sequence. The pink horizontal bars above the boxes and some 3' sequence of special bars and some 3' sequence. The pink horizontal bars above the boxes and some 3' sequence of special bars and some 3' sequence. The pink horizontal bars above the boxes and some 3' sequence. The pink horizontal bars above the boxes and some 3' sequence. The pink horizontal bars above the boxes and some 3' sequence. The pink horizontal bars above the boxes and some 3' sequence. The pink horizontal bars above the boxes and some 3' sequence. The pink horizontal bars above the boxes and some 3' sequence. The pink horizontal bars above the boxes and some 3' sequence. The pink horizontal bars above the boxes are the pink horizontal bars are the pink horizontal bars above the boxes are the pink horizontal bars above the boxes are the pink horizontal bars are the basis of the phenotypic overlap with Townes-Brocks, Okihiro and Feingold syndromes, we analyzed SALL1 (ref. 2), SALL4 (ref. 3) and MYCN4 but found no mutations in any of these genes (Supplementary Methods) of genomic DNA from the most severely affected individual (case 1, with lower lid coloboma, epilepsy and syringsomyelia) and identified a heterozygous deletion of 37.9-50.7 kb on Xq28, which removed exons 1 and 2 of FAM58A (Fig. 11,i). Using real-time PCR, we confirmed the deletion in the child and excluded it in her unaffected parents (Supplementary Methods and Supplementary Table 1 online). Through CGH with a customized oligonucleotide array enriched in probes for Xq28, followed by breakpoint cloning, we defined the exact deletion size as 40,068 bp (g.152,514,164_152,554,231del(chromosome X, NCBI Build 36.2); Fig. 1 jand Supplementary Figs. 2,3 online). The deletion removes the coding regions of exons 1 and 2 as well as intron 1 (2,774 bp), 492 bp of intron 2, and 36,608 bp of 5' sequence, including the 5' UTR and the entire KRT18P48 pseudogene (NCBI gene ID 340598). Paternity was proven using routine methods. We did not find deletions overlapping FAM58A in the available copy number variation (CNV) databases. Subsequently, we carried out qPCR analysis of the three other affected individuals (cases 2, 3 and 4) and the mother-daughter pair from the literature (cases 5 and 6). In case 3, we detected a de novo heterozygous felletion of 1,1–10.3 kb overlapping exon 5 (Supplementary Fig. 1b online). Using Xq28-targeted array CGH and breakpoint cloning, we identified a deletion of 4,249 by (g.152,504,371 del(chromosome X, NCBI Build 36.2); Fig. 1j and Supplementary Figs. 2,3), which removed 1,265 bp of intron 4, all of exon 5, including the 3' UTR, and 2,454 bp of 3' sequence. We found heterozygous FAM58A point mutations in the remaining cases (Fig. 1j. Supplementary Fig. 2, Supplementary Fig. 2, NCBI Build 36.2); Fig. 1j and Supplementary Table 1). In case 2, we identified the mutation 555+1G>A, affecting the splice donor site of intron 4. In case 3, we identified the mutation 500 heterozygous FAM58A point mutation 201dupT, which immediately results in a premature stop codon N68XfsX1. In cases 5 and 6, we detected the mutation 500 heterozygous FAM58A point mutations of PCR and sequencing or by qPCR. We confirmed paternity and de novo status of the point mutation 201dupT, which immediately results in a premature stop codon N68XfsX1. In case 3, we identified the mutation 500 heterozygous FAM58A point mutations of PCR and sequencing or by qPCR. We confirmed paternity and de novo status of the point mutation 201dupT, which immediately results in a premature stop codon N68XfsX1. In case 2, we identified the mutation 500 heterozygous FAM58A point mutation 500 heterozygous FAM58A point mutation 500 heterozygous FAM58A point mutations in the remaining example (Fig. 1). Supplementary Fig. 2, 3), which removed 1,265 bp of intron 4, all of example fig. 2, 3), which removed 1,265 bp of intron 4, all of example fig. 2,3), which removed 1,265 bp of intron 4, all of example fig. 2,3), which removed 1,265 bp of intron 4, all of example fig. 2,3), which are supplementary Fig. 2,3), which removed 1,265 bp of intron 4, all of example fig. 2,3), which are supplementary Fig. 2,3), which are suppl the mutations and deletions in all sporadic cases. None of the mutation swere seen in the DNA of 60 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were seen in the DNA of 60 unaffected female controls, and no larger deletions involving FAM58A were seen in the DNA of 60 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were seen in the DNA of 60 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female controls, and no larger deletions involving FAM58A were found in 93 unaffected female contro embryo and fetus. Genes homologous to FAM58A (NCBI HomoloGene: 13362) are found on the X chromosome in the chimpanzee and the dog. The zebrafish has a similar gene on chromosome 23. However, in the mouse and rat, there are no true homologs. These species have similar but intronless genes on chromosome omes 11 (mouse) and 10 (rat), most likely arising from a retrotransposon insertion event. On the murine X chromosome, the flanking genes Ato2b3 and Dusp9 are conserved, but only remnants of the FAM58A sequence can be detected. FAMS8A contains a cyclin-box-fold domain, a protein-binding domain found in cyclins with a role in cell cycle and transcription control. No human phenotype resulting from a cyclin pot receptor cells during embryogenesis?, 8. Cyclin D1 are viable but small and have reduced lifespan. They also have educed lifespan with a role in the retina, likely as a result of decreased cell proliferation and have reduced lifespan with a role in cell cycle and transcription control. No human phenotype resulting from a cyclin pot receptor cells during embryogenesis?, 8. Cyclin D1 are viable but small and have reduced lifespan with a role in cell cycle and transcription control. No human phenotype resulting from a cyclin pot receptor cells during embryogenesis?, 8. Cyclin D1 are viable but small and have reduced lifespan. They also have embryogenesis?, 8. Cyclin D1 are viable but small and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferation and have reduced lifespan with a result of decreased cell proliferati me developmental pathway. How do FAM58A mutations lead to STAR syndrome? Growth retardation (all cases; Table 1) and retinal abnormalities (three cases) are reminiscent of the reduced body size and retinal anomalies in cyclin D1 knockout mice?, 8. Therefore, a proliferation defect might be partly responsible for STAR syndrome. To address this question, we carried out a knockdown of FAM58A mRNA followed by a proliferation assay. Transfection of HEK293 cells with three different FAM58A-specific RNAi oligonucleotides for STAR syndrome. To address this question, we carried out a knockdown of FAM58A mRNA followed by a proliferation assay. Transfection of HEK293 cells with three different FAM58A-specific RNAi oligonucleotides for STAR syndrome. To address this question, we carried out a knockdown of FAM58A mRNA followed by a proliferation assay. Transfection of HEK293 cells with three different FAM58A-specific RNAi oligonucleotides for STAR syndrome. To address this question, we carried out a knockdown of FAM58A mRNA followed by a proliferation assay. Transfection of HEK293 cells with three different FAM58A-specific RNAi oligonucleotides for STAR syndrome. To address this question, we carried out a knockdown of FAM58A mRNA followed by a proliferation assay. Transfection of HEK293 cells with three different FAM58A-specific RNAi oligonucleotides for STAR syndrome. To address this question, search reminiscent of the reduced body size and retinal abnormalities (three cases) are reminiscent of the reduced body size and retinal abnormalities (three cases) are reminiscent of the reduced body size and retinal abnormalities (three cises) are reminiscent of the reduced body size and retinal abnormalities (three cises) are reminiscent of the reduced body size and retinal abnormalities (three cises) are reminiscent of the reduced body size and retinal abnormalities (three cises) are reminiscent of the reduced body size and retinal abnormalities (three cises) are reminiscent of the reduced body size and retinal abnormalities (thr egulation of FAM58A by MYCN. FAM58A is located approximately 0.56 Mb centromeric to MECP2 and FAM58A have been observed to date 13. Although other genes between FAM58A have been described and are not associated with a clinical phenotype in females 12, but no deletions overlapping both MECP2 and FAM58A have been observed to date 13. Although other genes between FAM58A have been described and are not associated with a clinical phenotype in females 12, but no deletions overlapping both MECP2 and FAM58A have been observed to date 13. Although other genes between FAM58A have been observed to date 13. Although other genes between FAM58A have been observed to date 13. Although other genes between FAM58A have been observed to date 13. Although other genes between FAM58A have been observed to date 13. Although other genes between FAM58A have been observed to date 13. Although other genes between FAM58A have been observed to date 13. Although other genes between FAM58A have been observed to date 13. Although other genes have been observed to date 13. A

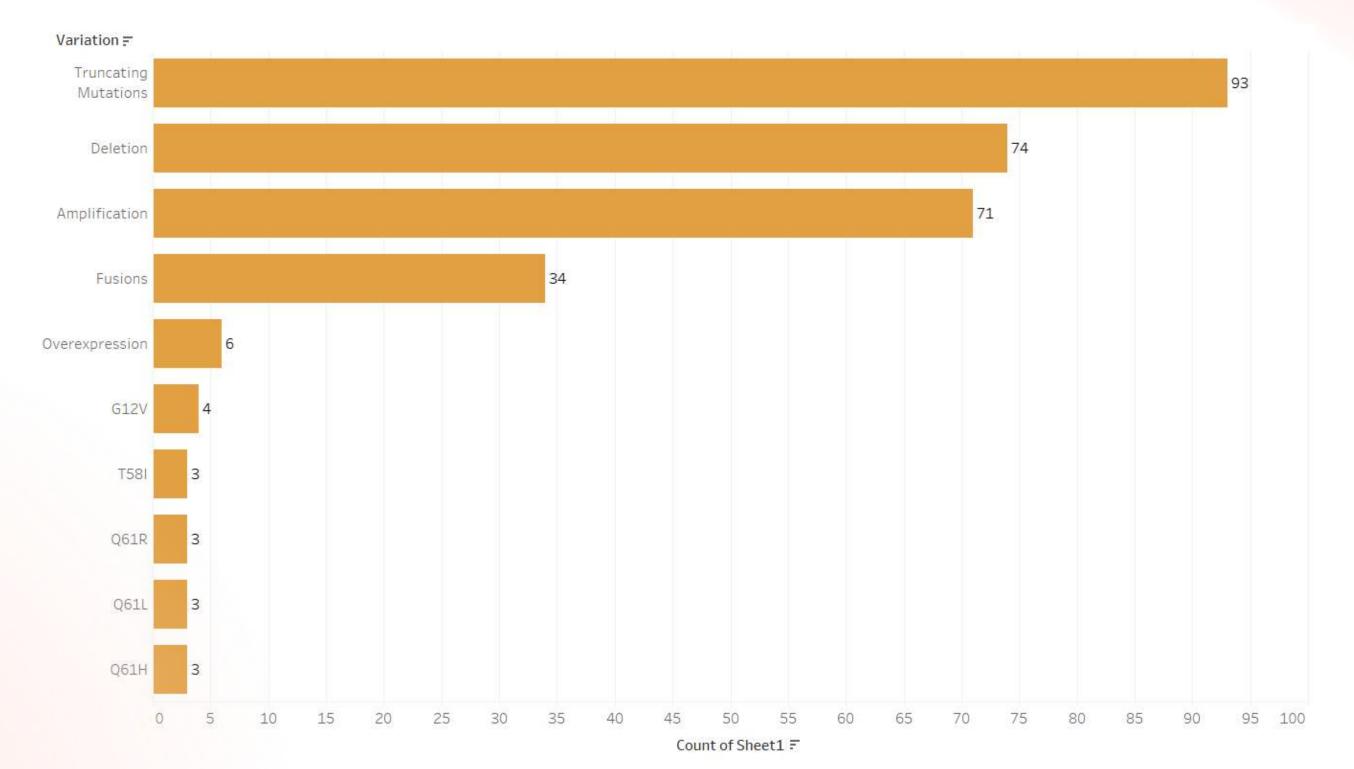
• จำนวนแต่ละ Class ในข้อมูล



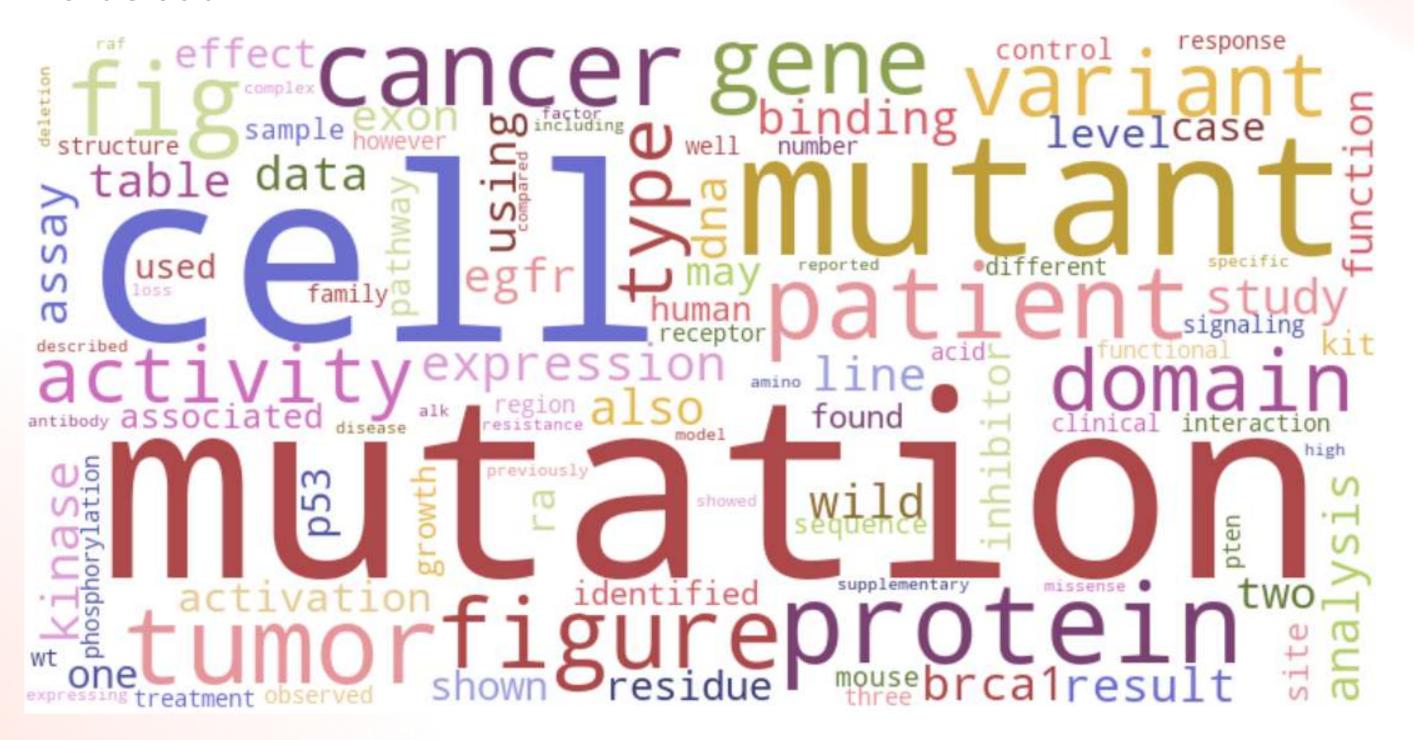
• จำนวน Gene ที่พบมาก 10 อันดับ



• จำนวน Variation ที่พบมาก 10 อันดับ



Wordcloud



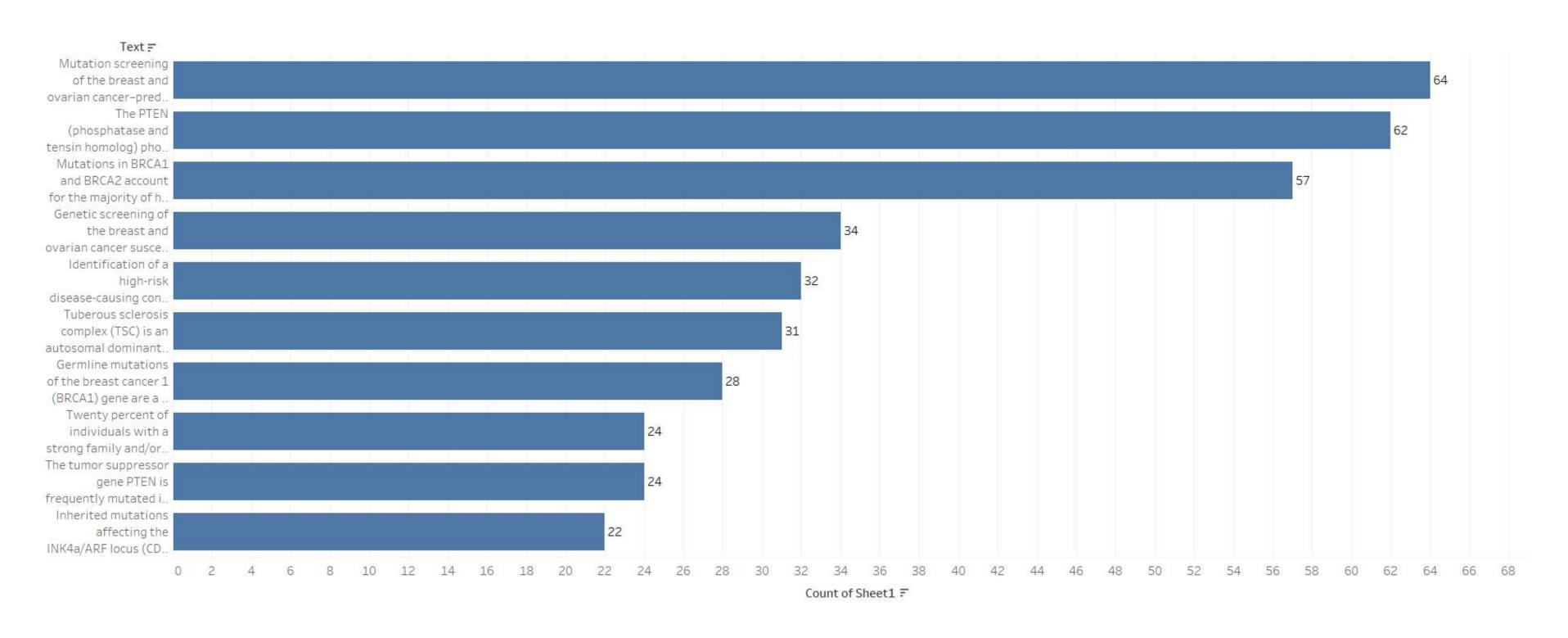
10 อันดับ

1 freq_text.head(10)

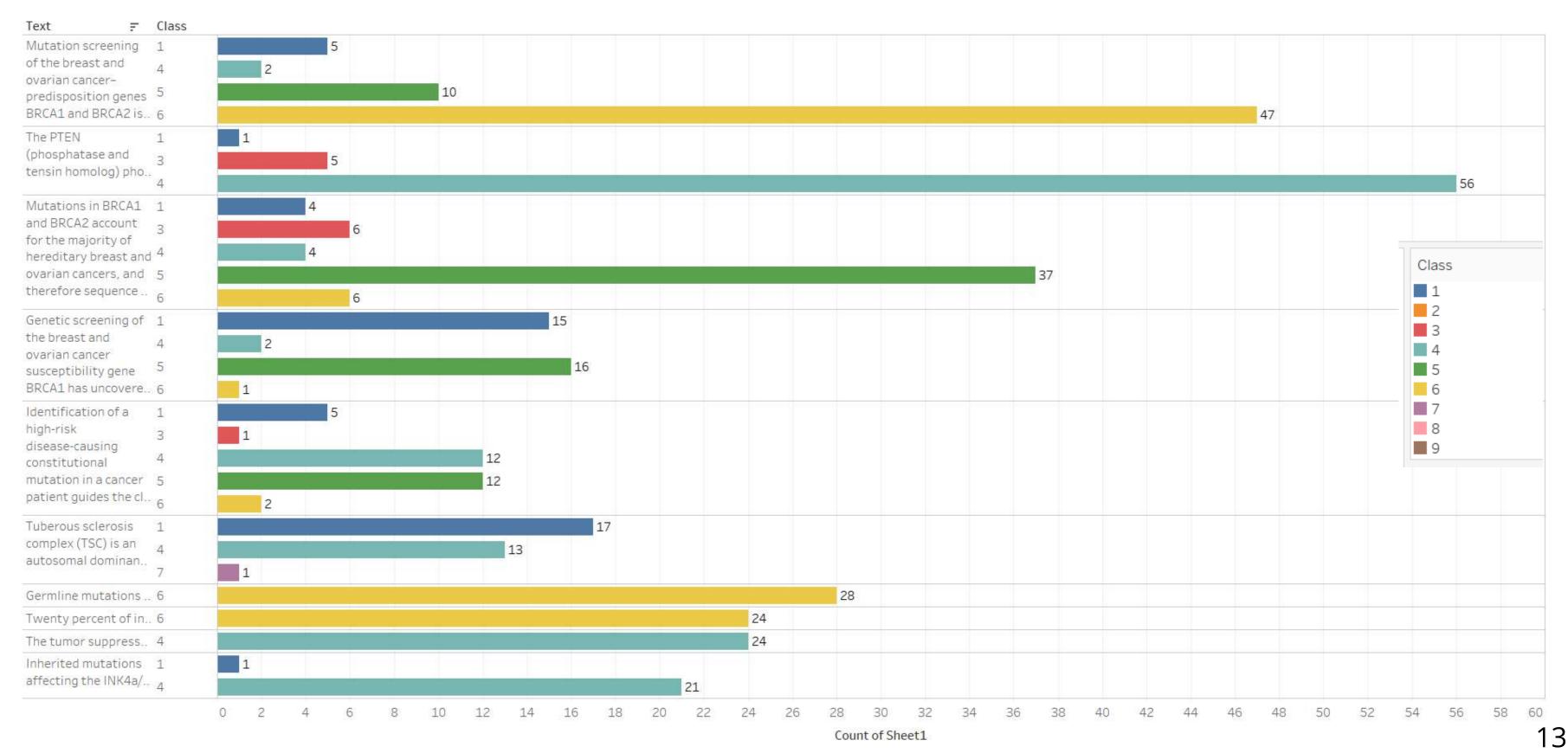
freq

token		
mutation	343839	
cell	312489	
mutant	126577	
protein	125359	
tumor	122770	
cancer	120094	
patient	111373	
fig	107416	
figure	104767	
gene	102268	

จำนวน Text ที่พบมาก 10 อันดับ



จำนวน Text ที่พบมาก 10 อันดับ ในแต่ละ Class



• มี Text ที่เหมือนกัน แต่อยู่คนละ Class

Text	Class	Gene	Variation	The PTEN	1	PTEN	G129A
Mutation screening	1	BRCA2	12627F	(phosphatase and	3	PTEN	V369G
of the breast and			G2748D	tensin homolog)			T131S
ovarian cancer-			E2663V	phosphatase is			K125R
predisposition genes			D3095E	unique in mammals in			1122V
BRCA1 and BRCA2 is			D2723G	terms of its tumor			A121V
becoming an	4	BRCA2	R2659K	suppressor activity,	4	PTEN	Y68D
increasingly			R2336H	exerted by			Y155C
important part of	5	BRCA2	V2908G	dephosphorylation of			V343L
clinical practice.			S1172L	the lipid second			V217D
Classification of rare			R2973C	messenger PIP3			T167A
nontruncating			P1819S	(phosphatidylinositol			T160I
sequence variants in			K513R	3,4,5-trisphosphate),			T131L
these genes is			K2950N	which activates the			T131I
problematic, because			K2729N	phosphoinositide			T131A
it is not known			12285V	•			R161G
whether these subtle			E462G	3-kinase/Akt/mTOR			R15K
			D2665G	(mammalian target			R130K
changes alter	6	BRCA2	Y3098H	of rapamycin)			R130G
function sufficiently			Y3092C	oncogenic pathway.			R130A
to predispose cells to			W2626C	Loss-of-function			P96Q
cancer development.			V894I	mutations in the			P95L
Using data from the			V30791	PTEN gene are			P169H
Myriad Genetic			V2969M	frequent in human			N48K

สรุปสิ่งที่ได้จากการทำ EDA

- จำนวนข้อมูลในแต่ละ class ของข้อมูล train
 มีไม่เท่ากัน
- ข้อมูลในบาง Class มีจำนวนน้อยเกินไป
- column text ในบาง row เหมือนกันทั้งหมด
- ความหมายของ class ไม่ชัดเจน

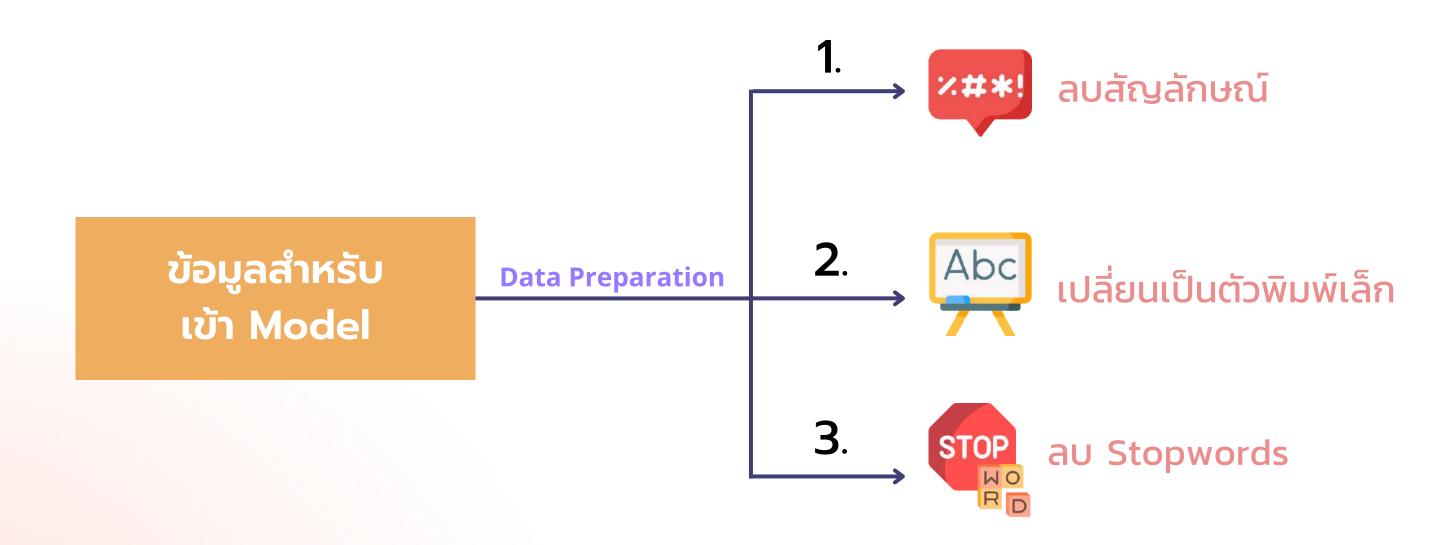


Objective

ใน Dataset นี้ แต่ละ Class (9 classes) สื่อถึงอะไร

🗕 โดยการทำ Topic Modeling

ี ลำดับกลุ่มที่ 2



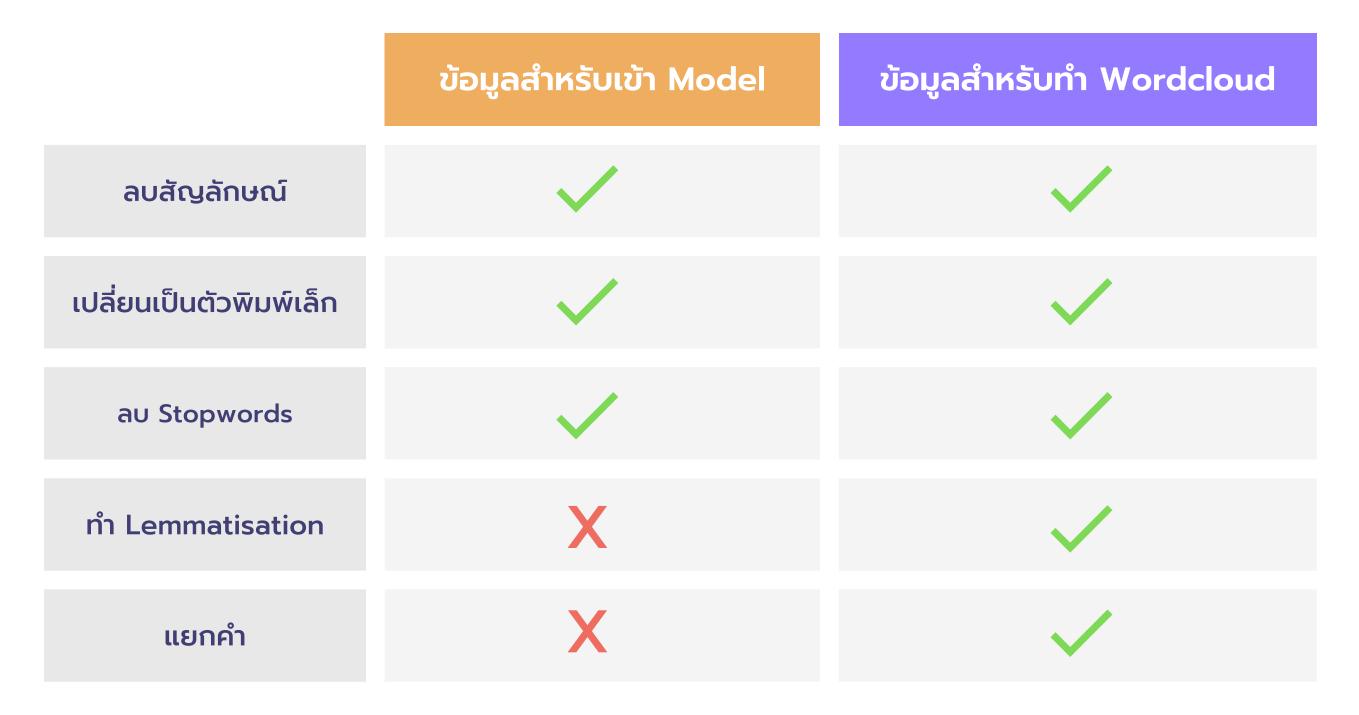
• ตัวอย่าง Stopwords

```
nltk.download('stopwords')
stopwords = set(nltk.corpus.stopwords.words('english'))
stopwords
[nltk data] Downloading package stopwords to /root/nltk data...
              Unzipping corpora/stopwords.zip.
[nltk data]
{'a',
                              'been',
 'about',
                              'before',
 'above',
                              'being',
 'after',
                              'below',
                              'between',
 'again',
                              'both',
 'against',
                              'but',
 'ain',
                              'by',
 'all',
                              'can',
 'am',
                              'couldn',
 'an',
                              "couldn't",
 'and',
                              'd',
 'any',
                              'did',
 'are',
                              'didn',
 'aren',
                              "didn't",
 "aren't",
                              'do',
 'as',
                              'does',
 'at',
                              'doesn',
                              "doesn't",
 'be',
                              'doing',
 'because',
```

• ข้อมูลที่ได้จาก Data Cleansing

3321 rows × 8 columns

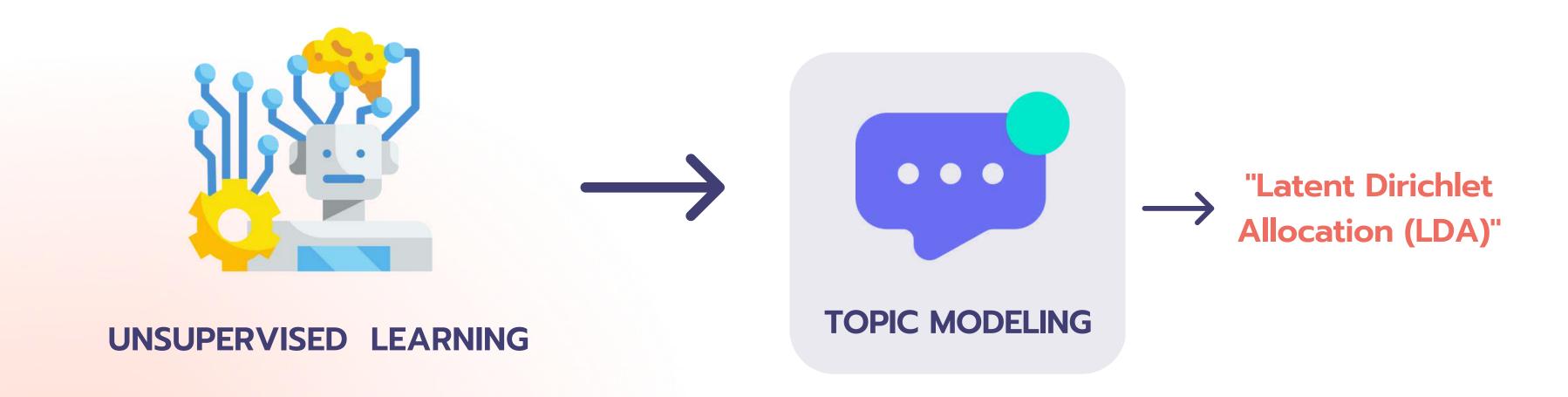
train_	_merge					
	ID	Gene	Variation	Class	TEXT	TEXT_cleaned
0	0	FAM58A	Truncating Mutations	1	Cyclin-dependent kinases (CDKs) regulate a var	cyclin dependent kinases cdks regulate variety
1	1	CBL	W802*	2	Abstract Background Non-small cell lung canc	abstract background non small cell lung cance
2	2	CBL	Q249E	2	Abstract Background Non-small cell lung canc	abstract background non small cell lung cance
3	3	CBL	N454D	3	Recent evidence has demonstrated that acquired	recent evidence demonstrated acquired uniparen
4	4	CBL	L399V	4	Oncogenic mutations in the monomeric Casitas B	oncogenic mutations monomeric casitas b lineag
3316	3316	RUNX1	D171N	4	Introduction Myelodysplastic syndromes (MDS)	introduction myelodysplastic syndromes mds het
3317	3317	RUNX1	A122*	1	Introduction Myelodysplastic syndromes (MDS)	introduction myelodysplastic syndromes mds het
3318	3318	RUNX1	Fusions	1	The Runt-related transcription factor 1 gene (runt related transcription factor 1 gene runx1
3319	3319	RUNX1	R80C	4	The RUNX1/AML1 gene is the most frequent targe	runx1 aml1 gene frequent target chromosomal tr
3320	3320	RUNX1	K83E	4	The most frequent mutations associated with le	frequent mutations associated leukemia recurre



ข้อมูลสำหรับทำ Wordcloud

train_	train_merge										
	ID	Gene	Variation	Class	TEXT	length	TEXT_cleaned	tokens_word	tokens_word_lemma		
0	0	FAM58A	Truncating Mutations	1	Cyclin-dependent kinases (CDKs) regulate a var	39672.0	cyclin dependent kinases cdks regulate variety	[cyclin, dependent, kinases, cdks, regulate, v	[cyclin, dependent, kinase, cdks, regulate, va		
1	1	CBL	W802*	2	Abstract Background Non-small cell lung canc	36691.0	abstract background non small cell lung cance	[abstract, background, non, small, cell, lung,	[abstract, background, non, small, cell, lung,		
2	2	CBL	Q249E	2	Abstract Background Non-small cell lung canc	36691.0	abstract background non small cell lung cance	[abstract, background, non, small, cell, lung,	[abstract, background, non, small, cell, lung,		
3	3	CBL	N454D	3	Recent evidence has demonstrated that acquired	36238.0	recent evidence demonstrated acquired uniparen	[recent, evidence, demonstrated, acquired, uni	[recent, evidence, demonstrated, acquired, uni		
4	4	CBL	L399V	4	Oncogenic mutations in the monomeric Casitas B	41308.0	oncogenic mutations monomeric casitas b lineag	[oncogenic, mutations, monomeric, casitas, b,	[oncogenic, mutation, monomeric, casitas, b, l		
•••						•••	•••	394.	···		
3316	3316	RUNX1	D171N	4	Introduction Myelodysplastic syndromes (MDS)	73895.0	introduction myelodysplastic syndromes mds het	[introduction, myelodysplastic, syndromes, mds	[introduction, myelodysplastic, syndrome, md,		
3317	3317	RUNX1	A122*	1	Introduction Myelodysplastic syndromes (MDS)	40127.0	introduction myelodysplastic syndromes mds het	[introduction, myelodysplastic, syndromes, mds	[introduction, myelodysplastic, syndrome, md,		
3318	3318	RUNX1	Fusions	1	The Runt-related transcription factor 1 gene (36200.0	runt related transcription factor 1 gene runx1	[runt, related, transcription, factor, 1, gene	[runt, related, transcription, factor, 1, gene		
3319	3319	RUNX1	R80C	4	The RUNX1/AML1 gene is the most frequent targe	32520.0	runx1 aml1 gene frequent target chromosomal tr	[runx1, aml1, gene, frequent, target, chromoso	[runx1, aml1, gene, frequent, target, chromoso		
3320	3320	RUNX1	K83E	4	The most frequent mutations associated with le	67136.0	frequent mutations associated leukemia recurre	[frequent, mutations, associated, leukemia, re	[frequent, mutation, associated, leukemia, rec		
3321 ro	ws×9 c	columns									

03 Modeling





```
Topic #0:

ure domain binding proteins mutant dna activity wild mutants residues

Topic #1:

patients exon cancer deletion cases dna family genes ure associated

Topic #2:

tumor tumors cancer ure genes loss lines supplementary dna human

Topic #3:

p53 dna binding wild domain mutant ure site mutants levels

Topic #4:

binding smad3 complex proteins tsc2 wild ure amino domain activity
```



Result of Topic Modeling

```
Topic #0:
gene expression cancer fusion genes ure tumor tumors protein cases
Topic #1:
kinase mutations mutant activity domain mutants protein activation type mutation
Topic #2:
braf ras raf met mutations activation kinase ure cancer activity
Topic #3:
patients patient mutations kinase clinical treatment domain fusion response gene
Topic #4:
resistance ure mutations mutation tumor pathway tumors treatment clinical identified
```



Result of Topic Modeling

```
Topic #0:
flt3 kinase dna sequence signaling activation expressing patients phosphorylation identified
Topic #1:
alk variants brcal activating variant patients functional kinase tscl assays
Topic #2:
brcal variants erbb4 assay ssa recombination repair variant missense mutants
Topic #3:
brcal variants assay vus functional neutral assays variant dna function
Topic #4:
mtor activating supplementary signaling amino mutants identified brcal variants wt
```



```
Topic #0:
dna patients cancer tumors genes binding domain site ure tumor
Topic #1:
activity catalytic residues kinase acid amino proteins mutant function phosphatase
Topic #2:
ure tumor genes cancer supplementary tumors mutant identified domain loss
Topic #3:
brca1 ras cancer activation activity assay ure function loss dna
Topic #4:
variants variant functional missense assays wild assay binding function activity
```



Result of Topic Modeling

```
Topic #0:
brca1 vus variants neutral deleterious brca2 functional assays assay sequence
Topic #1:
variants variant classification class pathogenic clinical family missense functional assay
Topic #2:
ure dna structure residues variants kinase binding supplementary proteins studies
Topic #3:
domain dna breast binding residues studies domains proteins number wild
Topic #4:
variants brca1 brct binding functional assays structure missense wild author
```



```
Topic #0:
cells ure dna erbb2 cell associated resistance activity binding expression
Topic #1:
brcal variants assay cells functional missense variant activity expression breast
Topic #2:
variants brcal deleterious brca2 vus family sequence variant neutral history
Topic #3:
patients expression cells tumor tumors ure cell clinical mutant 12
Topic #4:
binding site kinase resistance variants ure cells residues activation structure
```



Result of Topic Modeling

```
Topic #0:
egfr patients gefitinib lung response treatment cancer exon clinical tumor
Topic #1:
patients exon kit patient 11 clinical response study disease tumor
Topic #2:
egfr exon patients gefitinib response clinical treatment tumor survival study
Topic #3:
cancer tumor gene tumors genes ure dna pathway human binding
Topic #4:
ure mutants mutant signaling phosphorylation wt expressing receptor wild growth
```



Result of Topic Modeling

```
Topic #0:
idh1 idh2 2hg mutant idh alleles wild tumor mutation ure
Topic #1:
akt ure kinase activation mutant activity lines binding supplementary mutation
Topic #2:
mutant erbb2 genes tumors mutation kinase idh1 sf3b1 expressing associated
Topic #3:
erbb2 mutation kinase associated growth binding expressing mutant activity activation
Topic #4:
genes gene tumors methylation ure tumor observed rna supplementary lines
```



```
Topic #0:
sf3b1 genes splicing cell gene u2af35 exon mds supplementary rna
Topic #1:
ezh2 cell ure histone h3k27 levels methylation activity wt cancer
Topic #2:
sf3b1 splicing ure genes mds patients rna aberrant sf3b1mut cell
Topic #3:
ag cell alternative ezh2 methylation sf3b1 levels lines idh1 histone
Topic #4:
sf3b1 idh1 idh2 2hg splicing cell idh ure patients genes
```

05 Conclusion & Discussion



จากการทำ Topic Modeling พบว่า

- หลาย ๆ Class มีเนื้อหาที่ใกล้เคียงกัน
- แต่ในบาง Class มีเนื้อหาไปในทิศทางเดียวกัน เช่น Class 9 พูดถึง sf3b1 ชัดเจน
- ต้องอาศัยความรู้ในการตีความ Topic ค่อนข้างมาก

05 Conclusion & Discussion

ในการศึกษาครั้งหน้าจึงอยาก

- ลองใช้ grid search กับ LDA
- ลองนำข้อมูลที่ทำ Lemmatisation เข้า model ดู
- ลองใช้ Model อื่น ๆ ในการทำ Topic Modeling
- หาหน้าที่ และความเกี่ยวข้องกันของ Gene หรือ Variation ในแต่ละ Class



Responsible



Data Understanding
Modeling
Data interpretation



Data Understanding
Data Preparation
Data interpretation



Data Understanding Modeling Data interpretation

Thanks for your attention

Back Slides

02 Data Preparation

1. 💴 ลบสัญลักษณ์

```
[18] def tokenize(text):
    return re.sub('[^a-zA-Z0-9\n]', ' ', text) # เท่าที่เข้าใจคือจะแทนที่ทุกตัวที่ไม่ใช่ a-z, A-Z, 0-9, \n ด้วย ' '
[19] pipeline = [str, tokenize]

def prepare(text, pipeline):
    tokens = text
    for transform in pipeline:
        tokens = transform(tokens)
    return tokens

[20] train_merge['TEXT_cleaned'] = train_merge['TEXT'].apply(prepare, pipeline = pipeline) # ใช้ prepare() function
```

02 Data Preparation

2. [Abc] เปลี่ยนเป็นตัวพิมพ์เล็ก

```
[21] # ทำให้ str เป็นตัวพิมพ์เล็ก
train_merge['TEXT_cleaned'] = train_merge['TEXT_cleaned'].str.lower()
```

3. Stopwords

```
[22] # ลบ stopwords ออก
train_merge['TEXT_cleaned'] = train_merge['TEXT_cleaned'].str.split(' ').apply(lambda x: ' '.join(k for k in x if k not in stopwords))
```

02 Data Preparation

```
# ลบ spacebar หลาย ๆ ครั้งออก
     train_merge['TEXT_cleaned'] = train_merge['TEXT_cleaned'].str.replace(r'\s+', ' ', regex=True)
[24] # ตัดแบ่งคำ ไว้ใน col ใหม่ tokens_word
     train_merge['tokens_word'] = train_merge['TEXT_cleaned'].str.split()
```

ดูความยาวของ Text ในข้อมูล

```
1 # เพิ่ม col length
2 train_merge['length'] = train_merge['TEXT'].str.len()
```

length

39672.0

4 The RUNX1/AML1 gene is the most frequent targe... 32520.0

The most frequent mutations associated with le... 67136.0

	T	train_merge										
		ID	Gene	Variation	Class	TEXT						
	0	0	FAM58A	Truncating Mutations	1	Cyclin-dependent kinases (CDKs) regulate a var						
	1	1	CBL	W802*	2	Abstract Background Non-small cell lung canc						

R80C

K83E

36691.0 **CBL** Q249E Abstract Background Non-small cell lung canc... **CBL** Recent evidence has demonstrated that acquired... N454D 36238.0 L399V 4 Oncogenic mutations in the monomeric Casitas B... 41308.0 CBL D171N Introduction Myelodysplastic syndromes (MDS) ... RUNX1 **3317** 3317 RUNX1 A122* Introduction Myelodysplastic syndromes (MDS) ... 40127.0 The Runt-related transcription factor 1 gene (... 3318 RUNX1 Fusions 36200.0

3321 rows × 6 columns

RUNX1

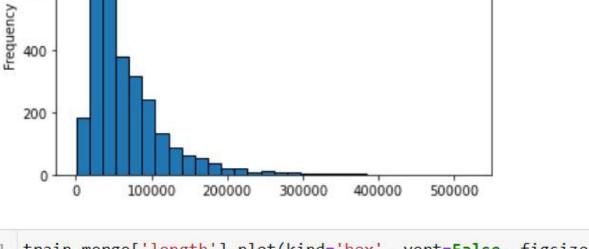
RUNX1

3319 3319

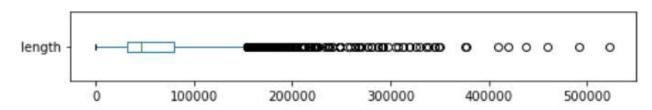
3320

3320

```
1 train_merge['length'].plot(kind='hist', bins=30, figsize=(6, 4), edgecolor='k')
<AxesSubplot:ylabel='Frequency'>
```







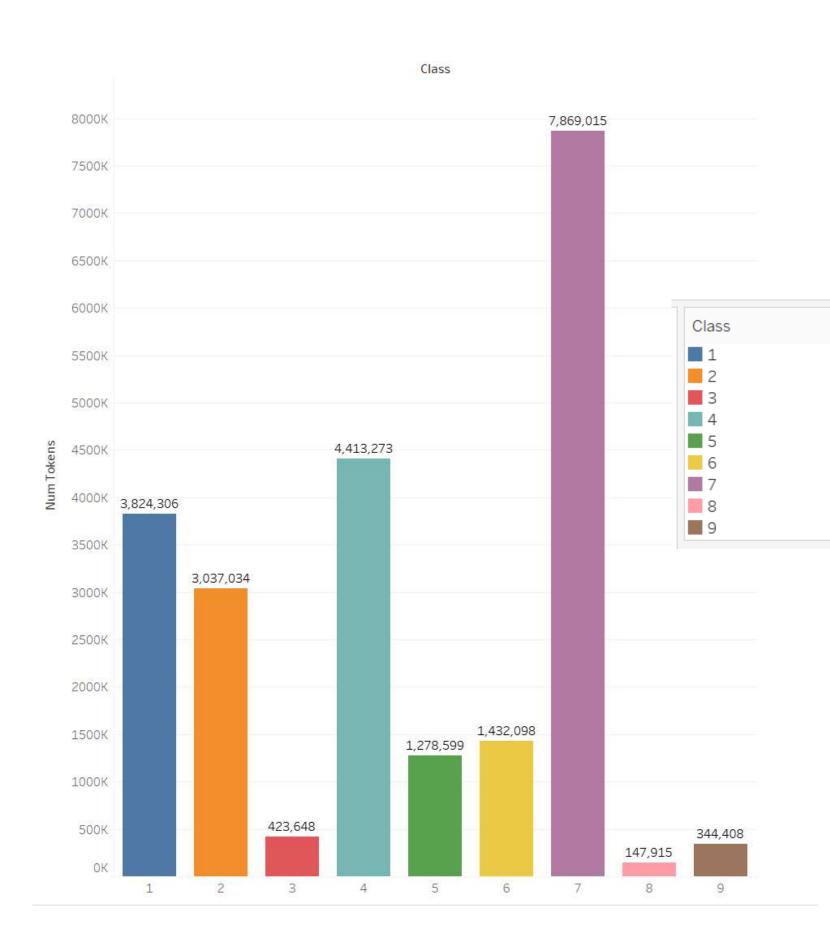
จำนวนToken ในแต่ละ Class

word len

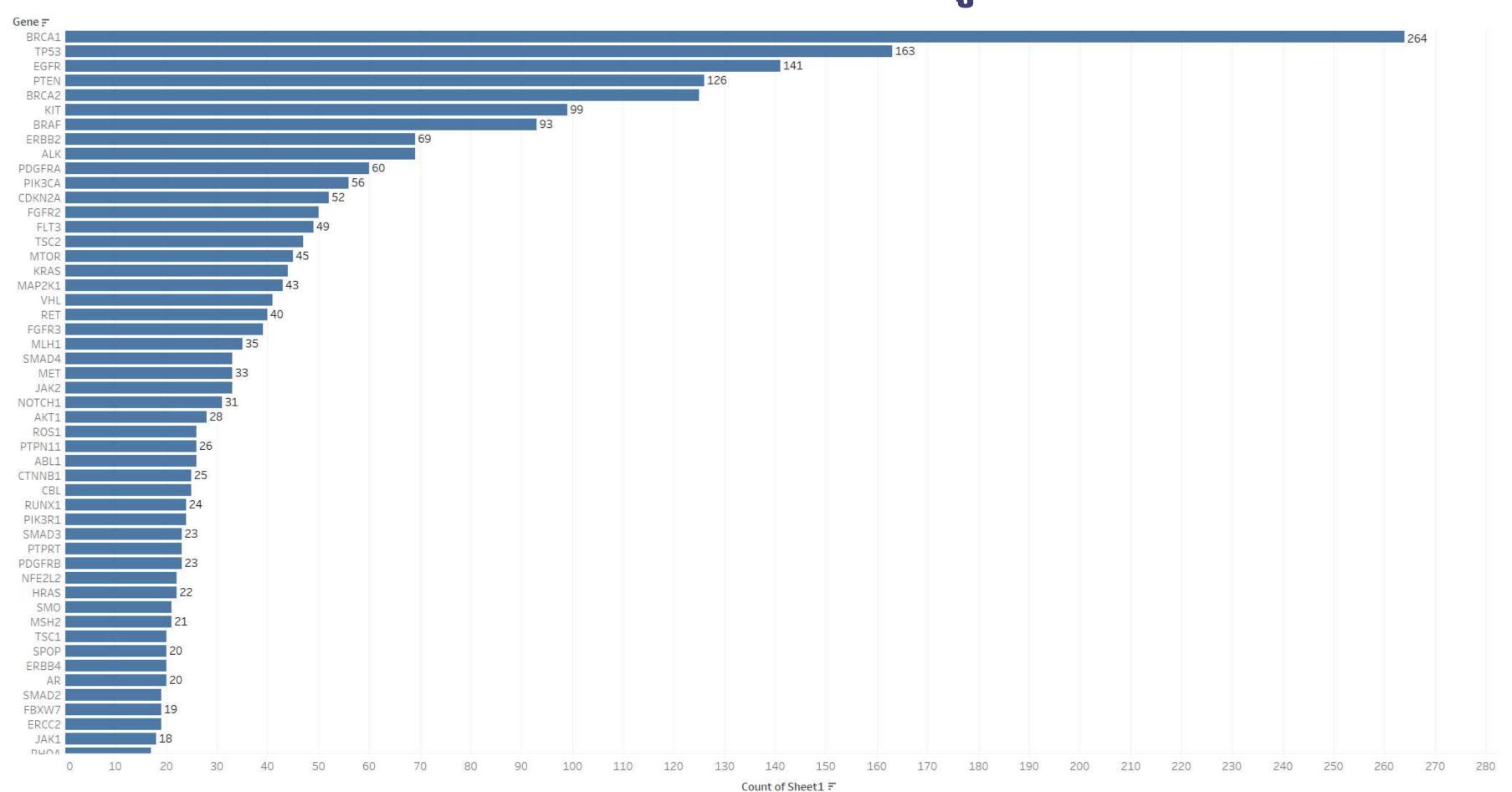
```
1 train_merge['num_tokens'] = train_merge['tokens_word_lemma'].map(len)
1 train_merge
```

99	ID	Gene	Variation	Class	TEXT	length	TEXT_cleaned	tokens_word	tokens_word_lemma	num_tokens
0	0	FAM58A	Truncating Mutations	1	Cyclin-dependent kinases (CDKs) regulate a var	39672.0	cyclin dependent kinases cdks regulate variety	[cyclin, dependent, kinases, cdks, regulate, V	[cyclin, dependent, kinase, cdks, regulate, va	4370
1	1	CBL	W802*	2	Abstract Background Non-small cell lung canc	36691.0	abstract background non small cell lung cance	[abstract, background, non, small, cell, lung,	[abstract, background, non, small, cell, lung,	4139
2	2	CBL	Q249E	2	Abstract Background Non-small cell lung canc	36691.0	abstract background non small cell lung cance	[abstract, background, non, small, cell, lung,	[abstract, background, non, small, cell, lung,	4139
3	3	CBL	N454D	3	Recent evidence has demonstrated that acquired	36238.0	recent evidence demonstrated acquired uniparen	[recent, evidence, demonstrated, acquired, uni	[recent, evidence, demonstrated, acquired, uni	3841
4	4	CBL	L399V	4	Oncogenic mutations in the monomeric Casitas B	41308.0	oncogenic mutations monomeric casitas b lineag	[oncogenic, mutations, monomeric, casitas, b,	[oncogenic, mutation, monomeric, casitas, b, l	4254
	444	344	***	***	100	2944	***			
3316	3316	RUNX1	D171N	4	Introduction Myelodysplastic syndromes (MDS)	73895.0	introduction myelodysplastic syndromes mds het	[introduction, myelodysplastic, syndromes, mds	[introduction, myelodysplastic, syndrome, md,	8153
3317	3317	RUNX1	A122*	1	Introduction Myelodysplastic syndromes (MDS)	40127.0	introduction myelodysplastic syndromes mds het	[introduction, myelodysplastic, syndromes, mds	[introduction, myelodysplastic, syndrome, md,	4495
3318	3318	RUNX1	Fusions	1	The Runt-related transcription factor 1 gene (36200.0	runt related transcription factor 1 gene runx1	[runt, related, transcription, factor, 1, gene	[runt, related, transcription, factor, 1, gene	4593
3319	3319	RUNX1	R80C	4	The RUNX1/AML1 gene is the most frequent targe	32520.0	runx1 aml1 gene frequent target chromosomal tr	[runx1, aml1, gene, frequent, target, chromoso	[runx1, aml1, gene, frequent, target, chromoso	3465
3320	3320	RUNX1	K83E	4	The most frequent mutations associated with le	67136.0	frequent mutations associated leukemia recurre	[frequent, mutations, associated, leukemia, re	[frequent, mutation, associated, leukemia, rec	7013

3321 rows × 10 columns



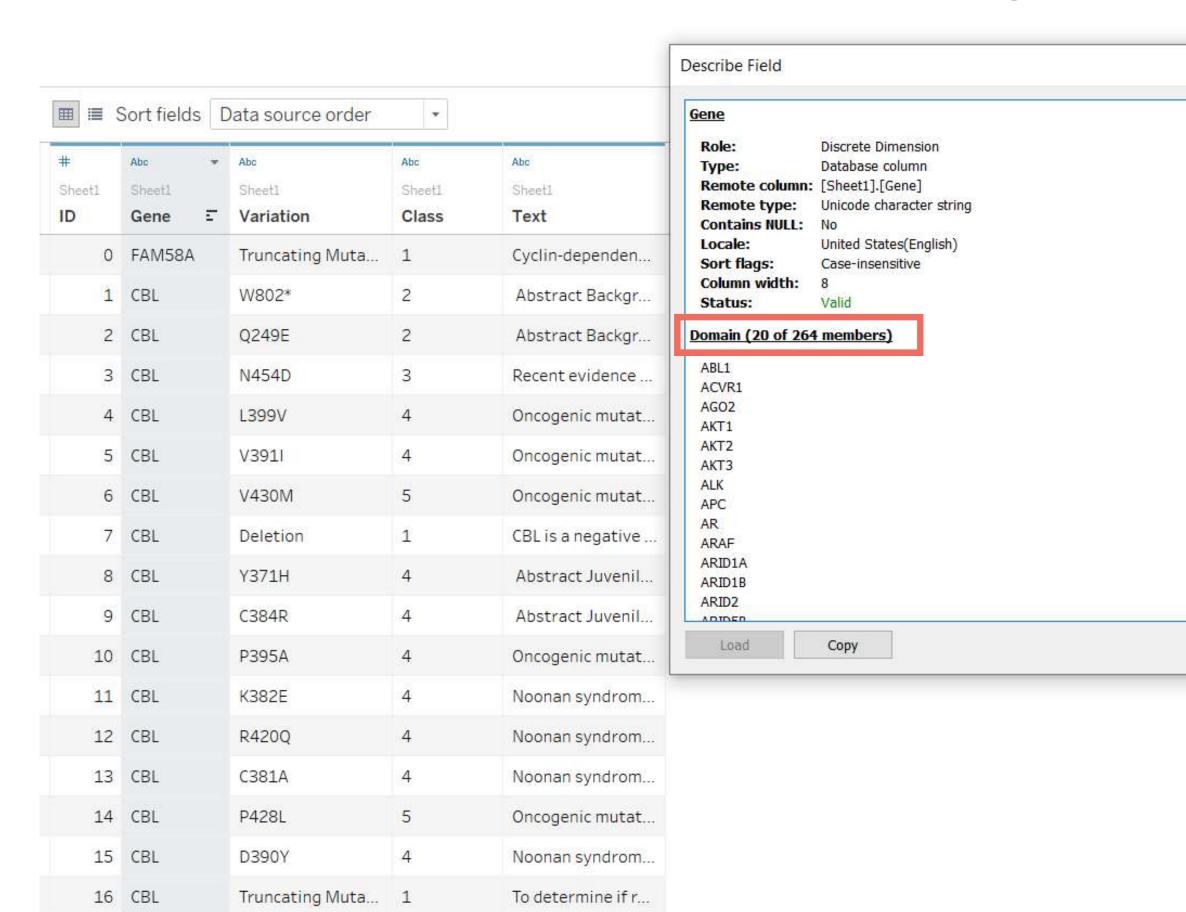
จำนวน Gene ที่พบในข้อมูล



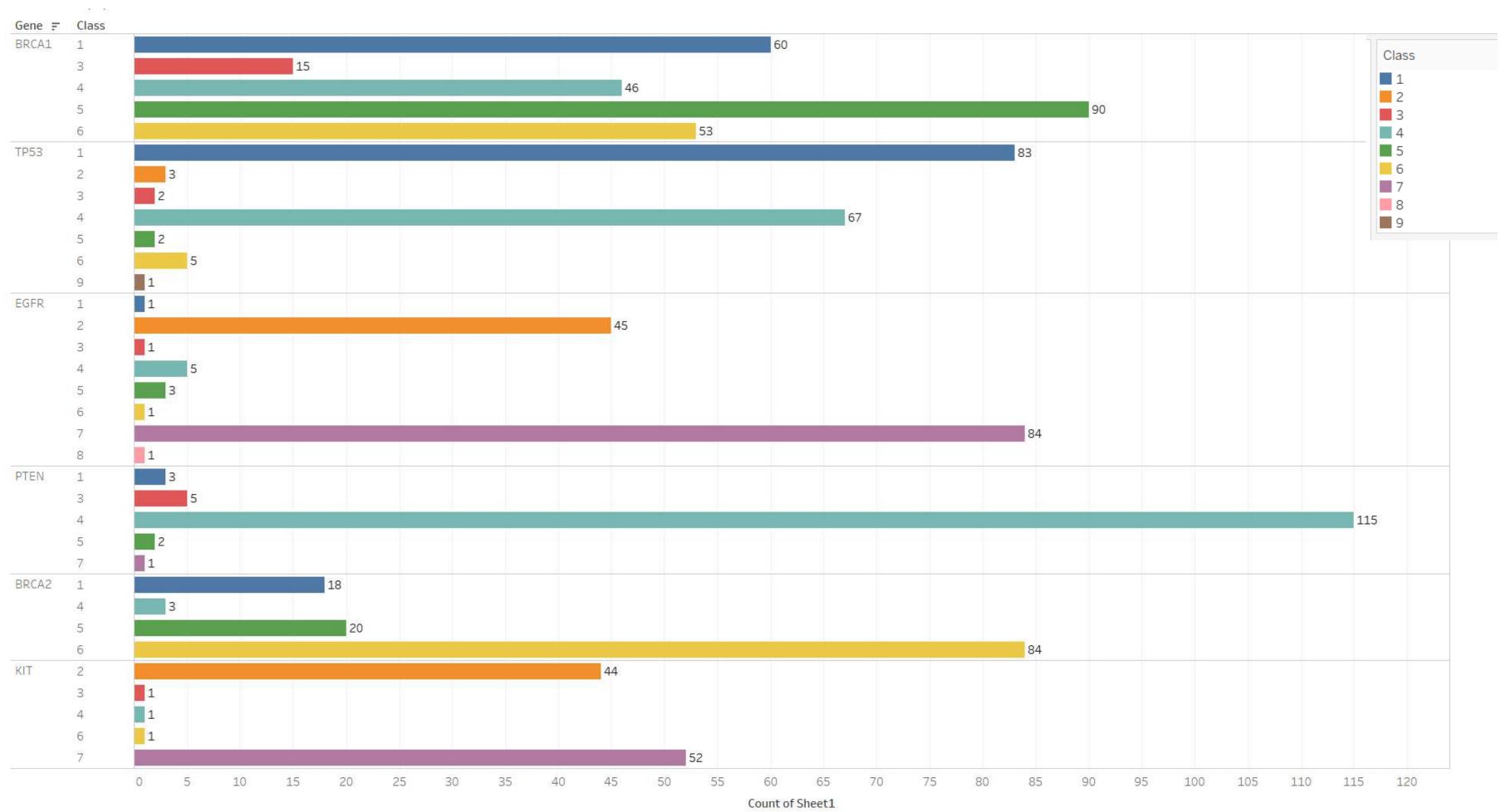
จำนวน Gene ที่พบในข้อมูล

×

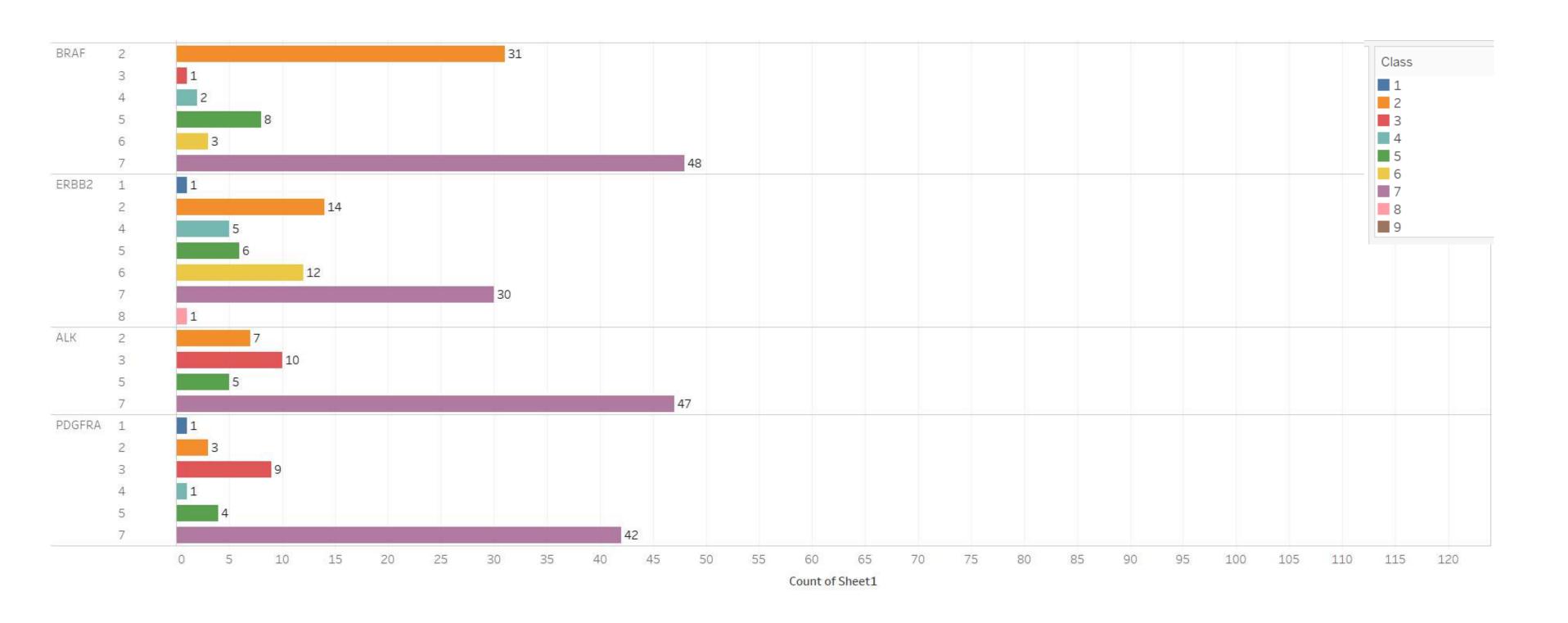
Close



จำนวน Gene ที่พบมาก 10 อันดับ ในแต่ละ Class



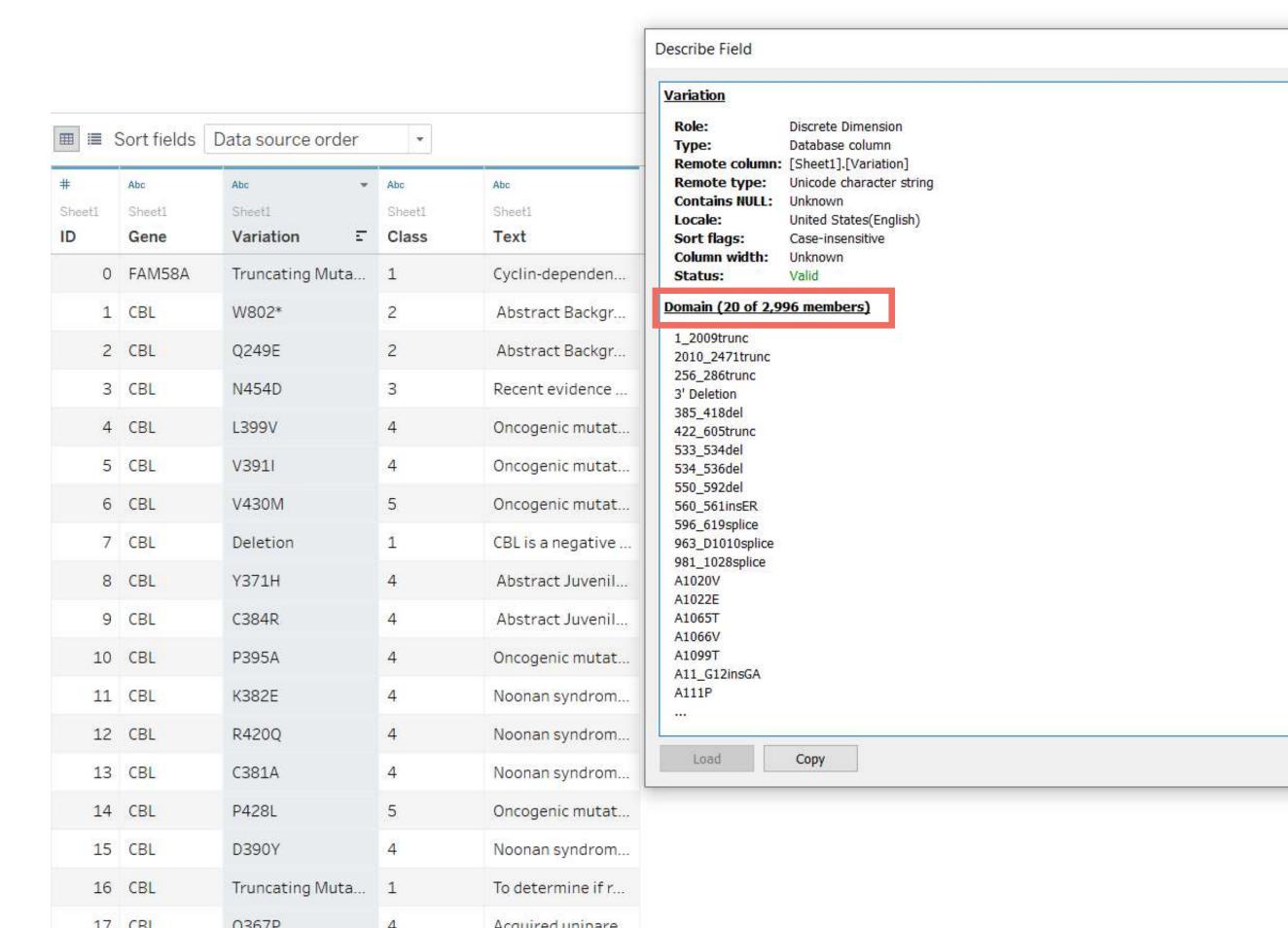
จำนวน Gene ที่พบมาก 10 อันดับ ในแต่ละ Class (ต่อ)



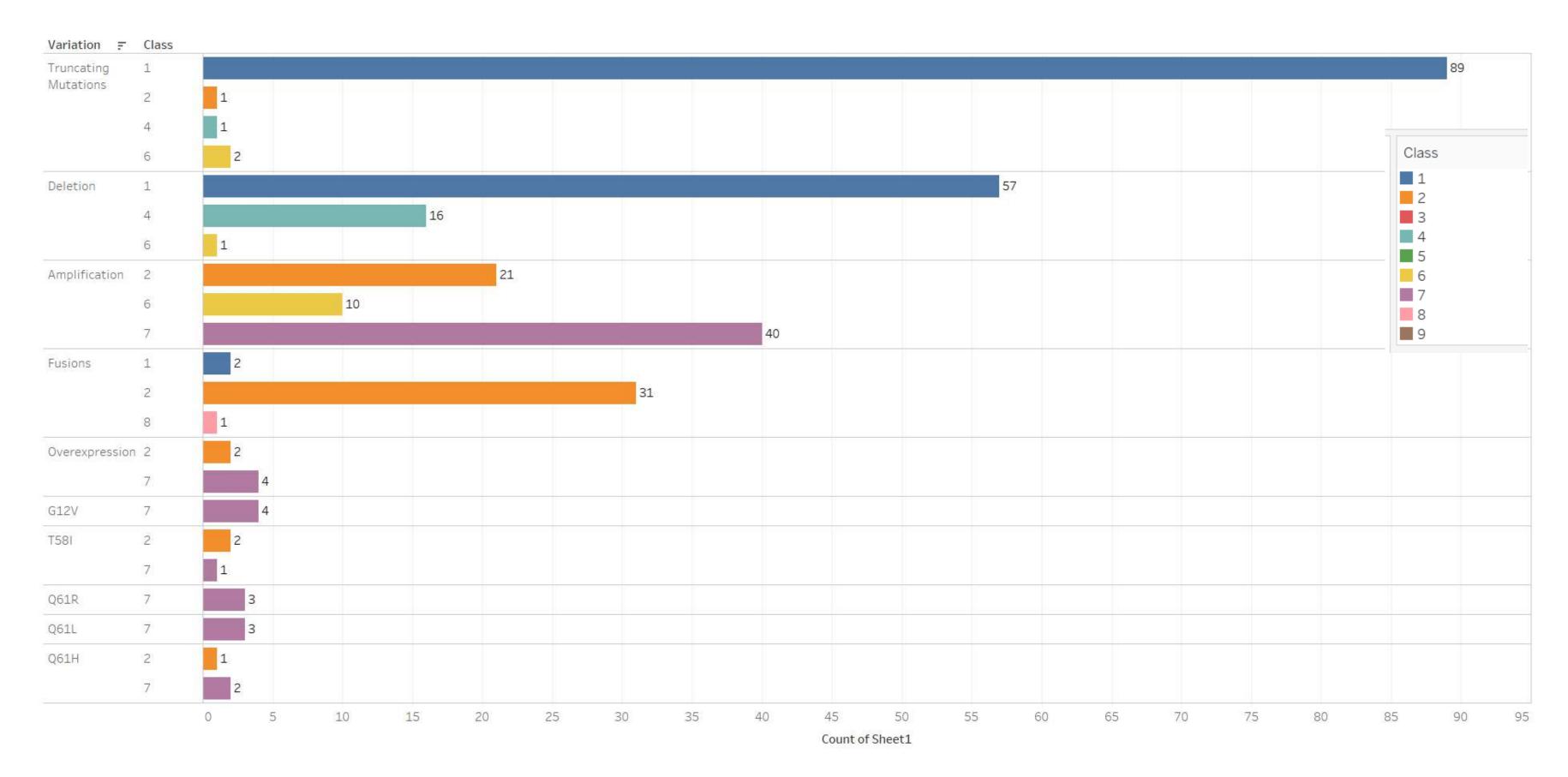
จำนวน Variation ที่พบในข้อมูล

X

Close



จำนวน Variation ที่พบมาก 10 อันดับ ในแต่ละ Class



จำนวน Text ที่พบในข้อมูล

