

TOPICS

Next

- AUC - ROC

- PR-curve

- KNN (cont)

- Imbalanced data

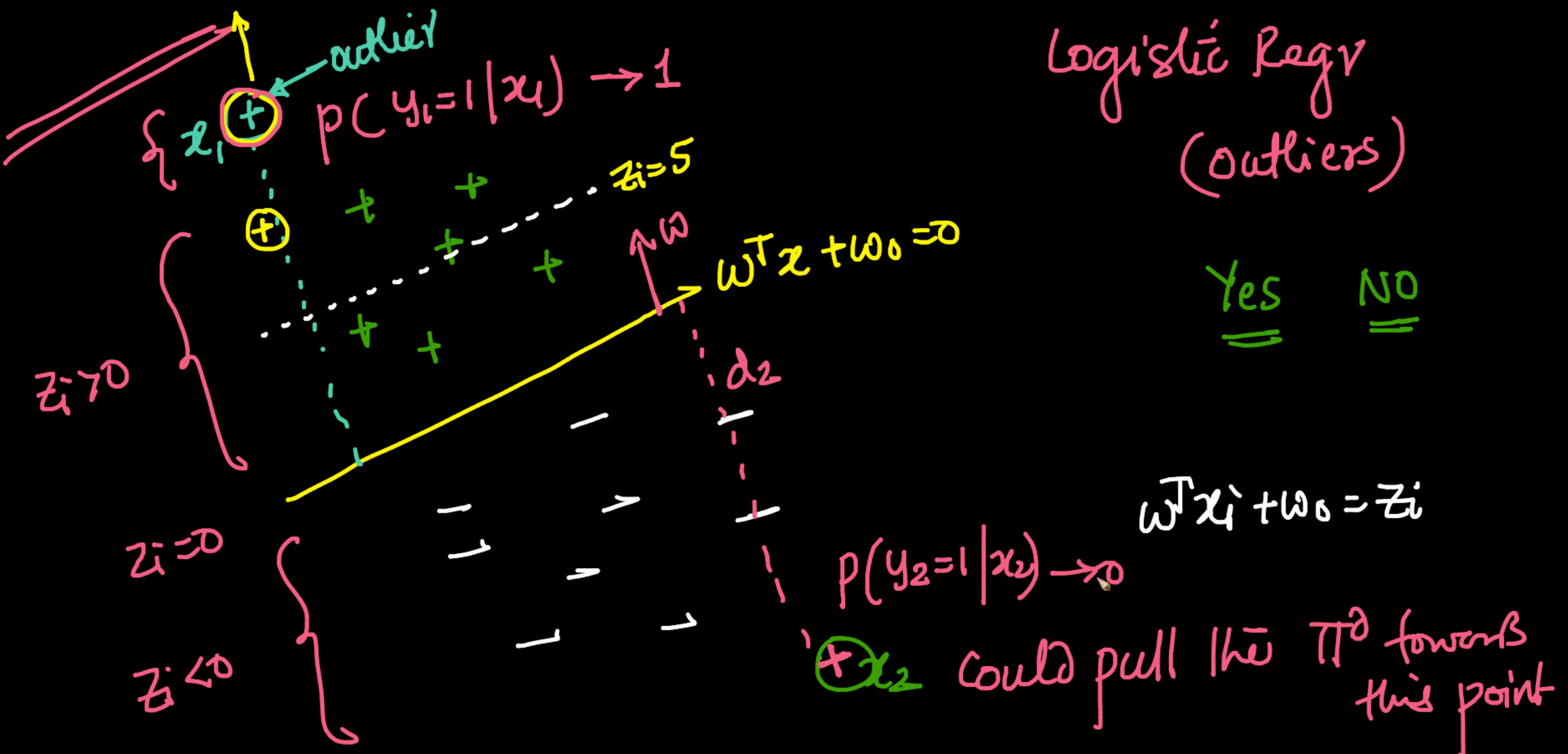
+ Odds ratio

KNN (in detail)

simple (intuitive)

→ outliers (visual)
Intuition + Math
Logistic regression (contd)
churn case-study (fix) → imbalanced-data
confusion Matrix ; Precision ; Recall ; F1

weights



Case 1: +ve outlier on the correct side

$w^T x_i + w_0$: large +ve

squashing: $\sigma(w^T x_i + w_0)$: close to +1

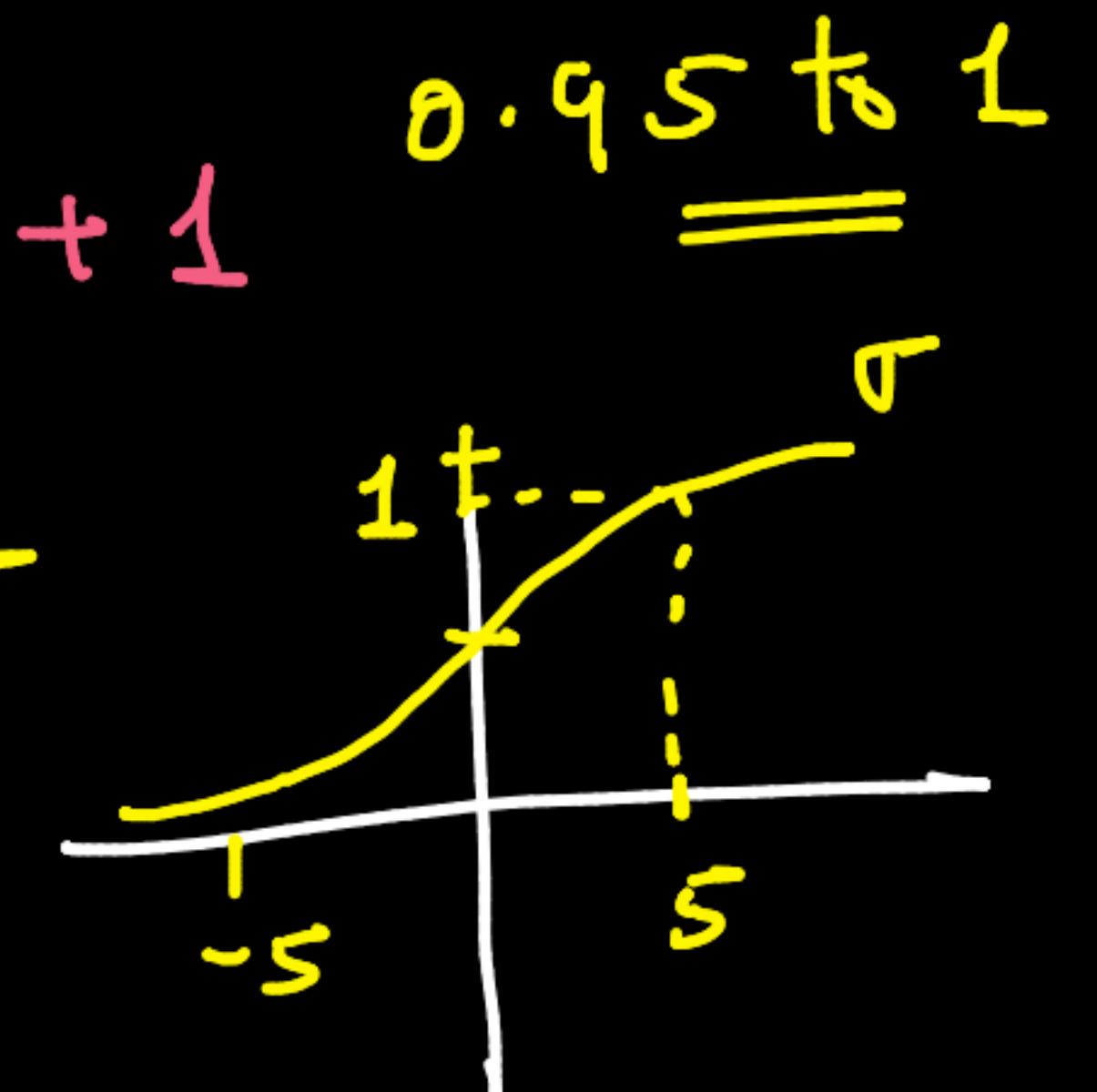
$y_i = 1$

↳ not changing a lot

{ Not much impact on π_{θ}^{sep}

log-loss: $-\log(\hat{y}_i)$

↳ ≈ 0



$$\text{opt: } \min_{\mathbf{w}} \sum_{i=1}^n \text{log-loss}_i + \lambda \|\mathbf{w}\|_2^2$$

Case 2: +ve labelled point ($y_i = 1$)
which is on wrong side

$$z_2 = \mathbf{w}^T \mathbf{x}_2 + w_0 : \text{large -ve value}$$

$$\sigma(z_2) : \text{close to zero} = \hat{y}_i$$

$$\text{log-loss}_2 = -\log(\hat{y}_i) : \text{large +ve value} =$$

Google

plot(-log(x))



All

Maps

Images

Videos

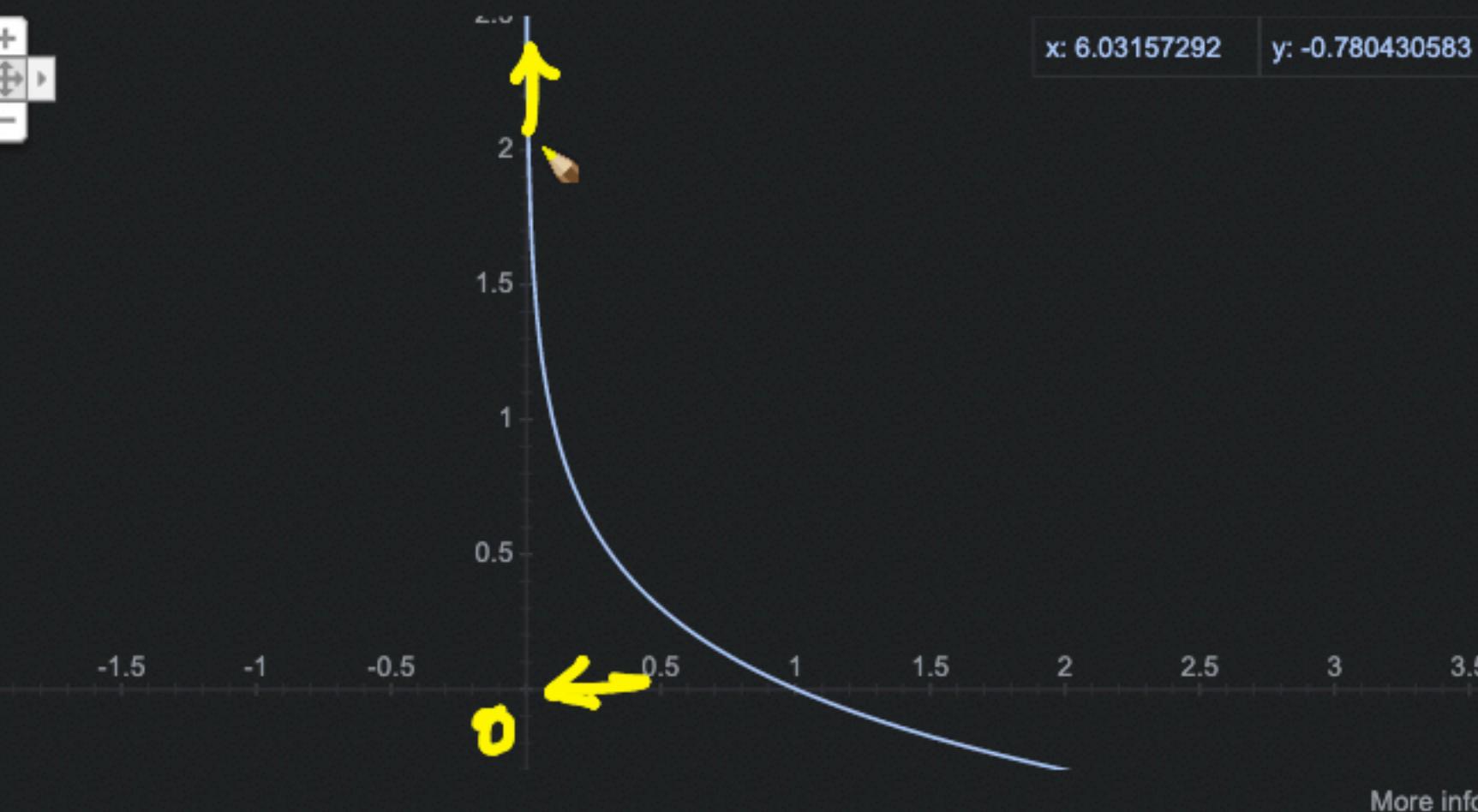
News

More

Tools



About 1,42,00,00,000 results (0.53 seconds)

Graph for $-\log(x)$ 

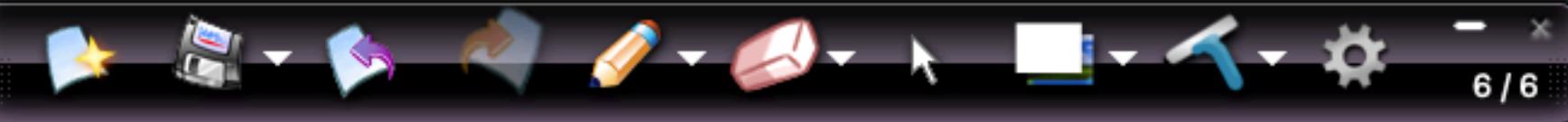
$$\hat{y}_i = 1$$
$$\hat{y}_i \rightarrow 0$$

$$\text{log-loss}_i = -\log(\hat{y}_i)$$

<https://www.rapidtables.com/math/algebra/logarithms.htm>**logarithm graph | graph of log(x) - RapidTables.com**

Graph of $\log(x)$. $\log(x)$ function graph. Logarithm graph. $y = f(x) = \log_{10}(x)$. $\log(x)$ graph properties. $\log(x)$ is defined for positive values of x . $\log(x)$...

Videos

Draw the graph of ' $y=-\log(x)$ ' when the graph of ' $y=\log(x)$ ' is ...YouTube · DoubtNut
31-May-2020

Opt: $\min_{w_j} \sum_{i=1}^n \log loss_i + \lambda \text{reg}$

w_j

$\log loss_2$ (outlier)

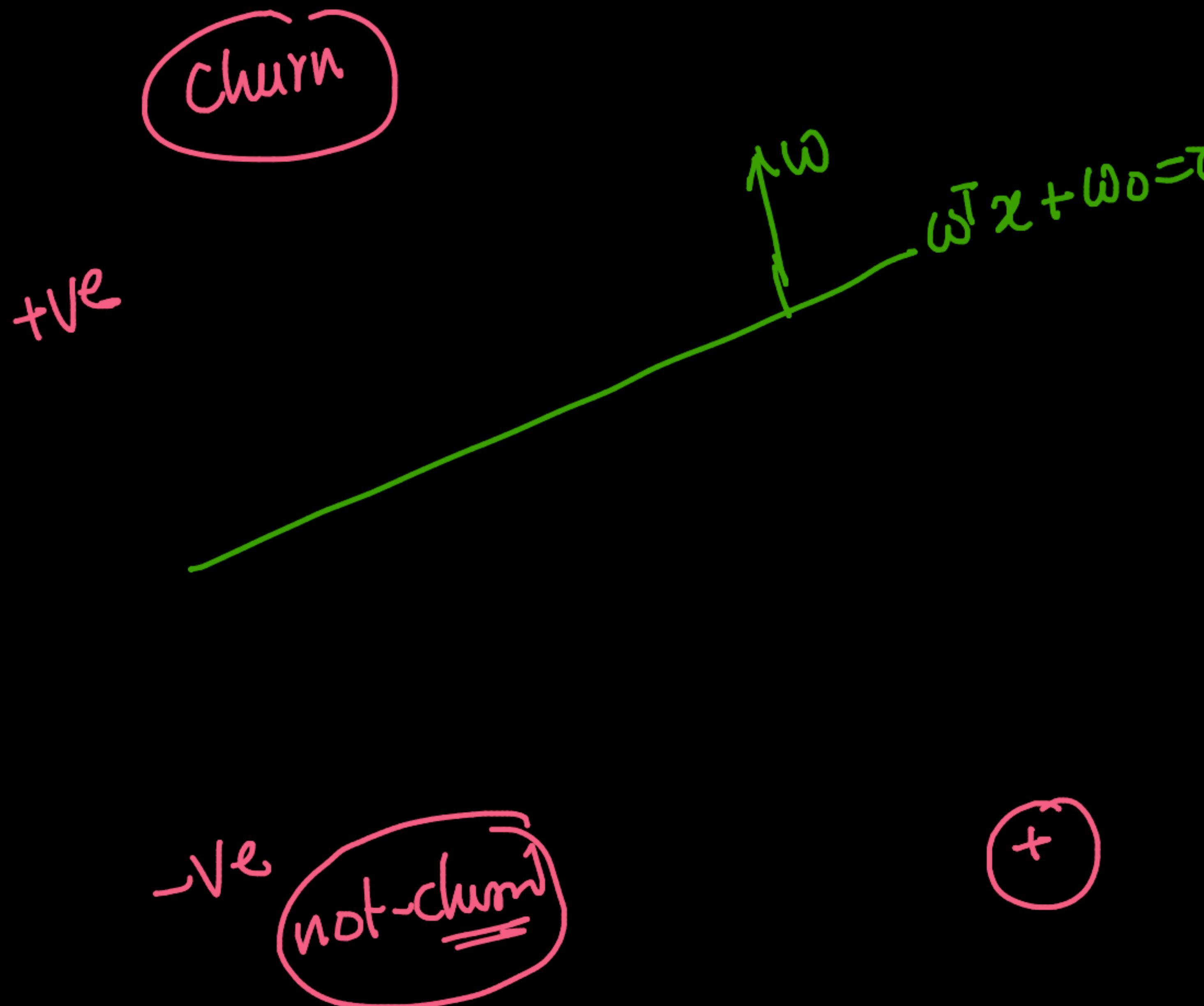
large true-value

fussle: x_2 & all other points

Savitz → some impact by far-away points
(outliers) on the wrong side of \hat{S}_{PT}

→ -log(\hat{y}_i)

Some impact
of such outliers
=



outliers,
→ scenarios
→ data issue

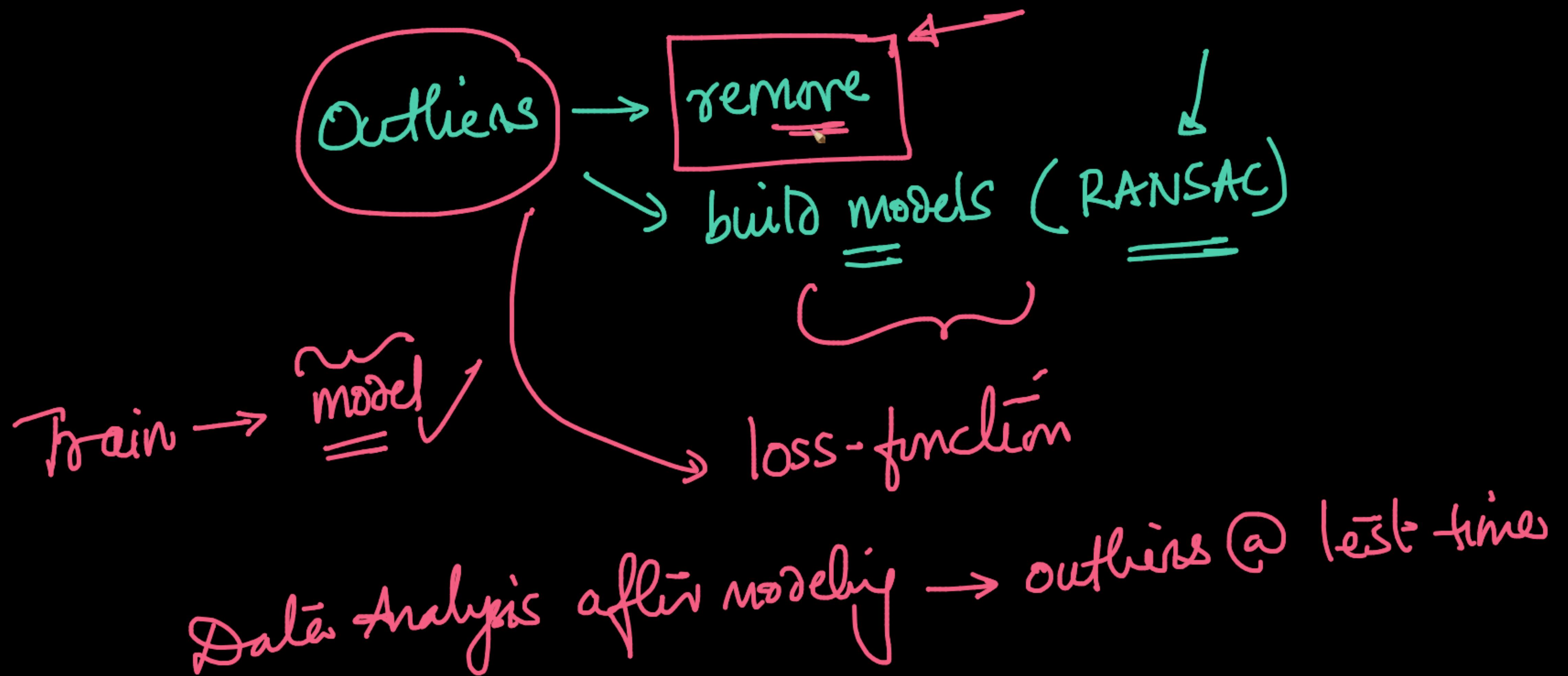
x_1 : died



no-sigmoid \rightarrow impact of outliers is significant

Sigmoid \rightarrow some outliers will not effect much

not all
all



Churn - problem

class 0 → 2850

[imbalanced - data]

Class 1 → 483

~ 85.5% of data class - 0

Q1

D-Test

A hand-drawn diagram of a neural network with three layers:

- Input Layer:** 2 neurons labeled x_1 and x_2 .
- Hidden Layer:** 3 neurons labeled y_1 , y_2 , and y_3 .
- Output Layer:** 2 neurons labeled \hat{y}_1 and \hat{y}_2 .

The diagram shows connections between neurons. A box labeled "loop" is connected to the first neuron of the hidden layer (y_1). The connections and their weights are as follows:

From	To	Weight
x_1	y_1	1
x_1	y_2	1
x_1	y_3	0
x_2	y_1	0
x_2	y_2	0
x_2	y_3	0
y_1	\hat{y}_1	1
y_1	\hat{y}_2	0
y_2	\hat{y}_1	0
y_2	\hat{y}_2	1
y_3	\hat{y}_1	0
y_3	\hat{y}_2	0

$$\hat{y}_i = \text{logistík füg}(x_i)$$

f_i

$\hat{y}_i = 0$ [bad model]

incorrect

accuracy: $\approx 85.5\%$

$\frac{85}{100}$

imbalanced \Rightarrow accuracy is a bad
metric

+ Code + Text

```
[ ] Pipeline(steps=[('standardscaler', StandardScaler()), ('logisticregression', LogisticRegression(C=0.0142836737608913))])
```

```
[ ] test_score = scaled_lr.score(X_test, y_test)  
print(test_score)
```

```
y_pred = scaled_lr.predict(X_test)
```

0.856071964017991

```
[ ] from sklearn.metrics import accuracy_score  
  
print(f"Accuracy : {accuracy_score(y_test, y_pred)*100}%)
```

Accuracy : 85.6071964017991%

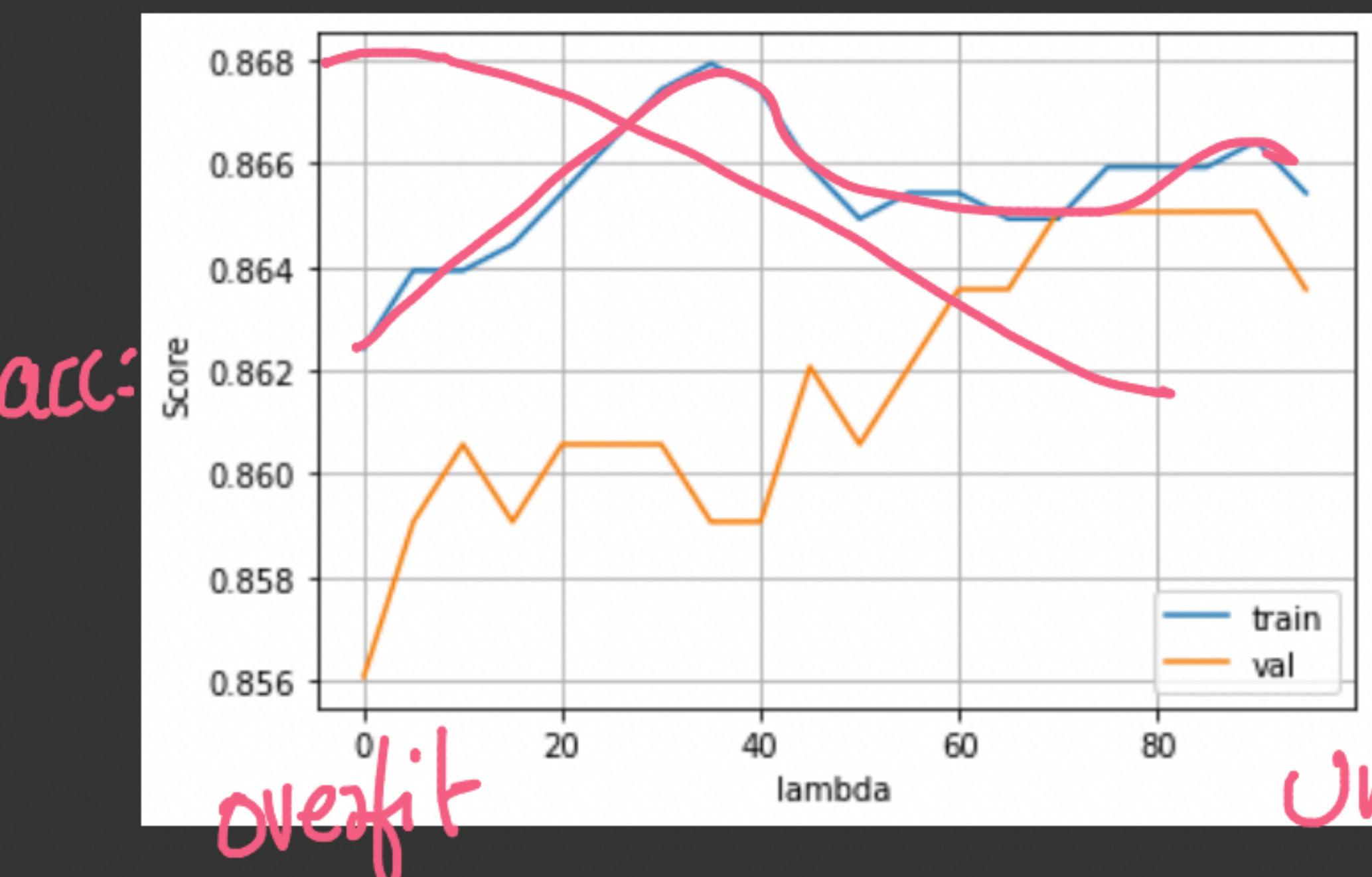
```
[ ] from sklearn.metrics import confusion_matrix  
  
confusion_matrix(y_test, y_pred)  
  
array([[560,   6],  
       [ 90,  11]])
```

↓
85.60% as accuracy

→ Same as a terrible
Model $\hat{y}_i = 0 \text{ or } 1$

+ Code + Text

[] plt.show()



[] np.argmax(val_scores)

14

[] val_scores[0]

0.856071964017991

Logistic Regression.ipynb - Colab | sklearn.linear_model.LogisticR... | Classifier comparison — scikit-learn 0.24.2 | plot(-log(x)) - Google Search | +

colab.research.google.com/drive/1G_D39wa40kVahULehDS4lWhB_PXhP-xJ#scrollTo=X2vZAFjWc-8E

+ Code + Text

[] plt.show()

Score

lambda

train

val

0.868

0.866

0.864

0.862

0.860

0.858

0.856

0 20 40 60 80

WRONG

{x}

acc on imbalanced data

14

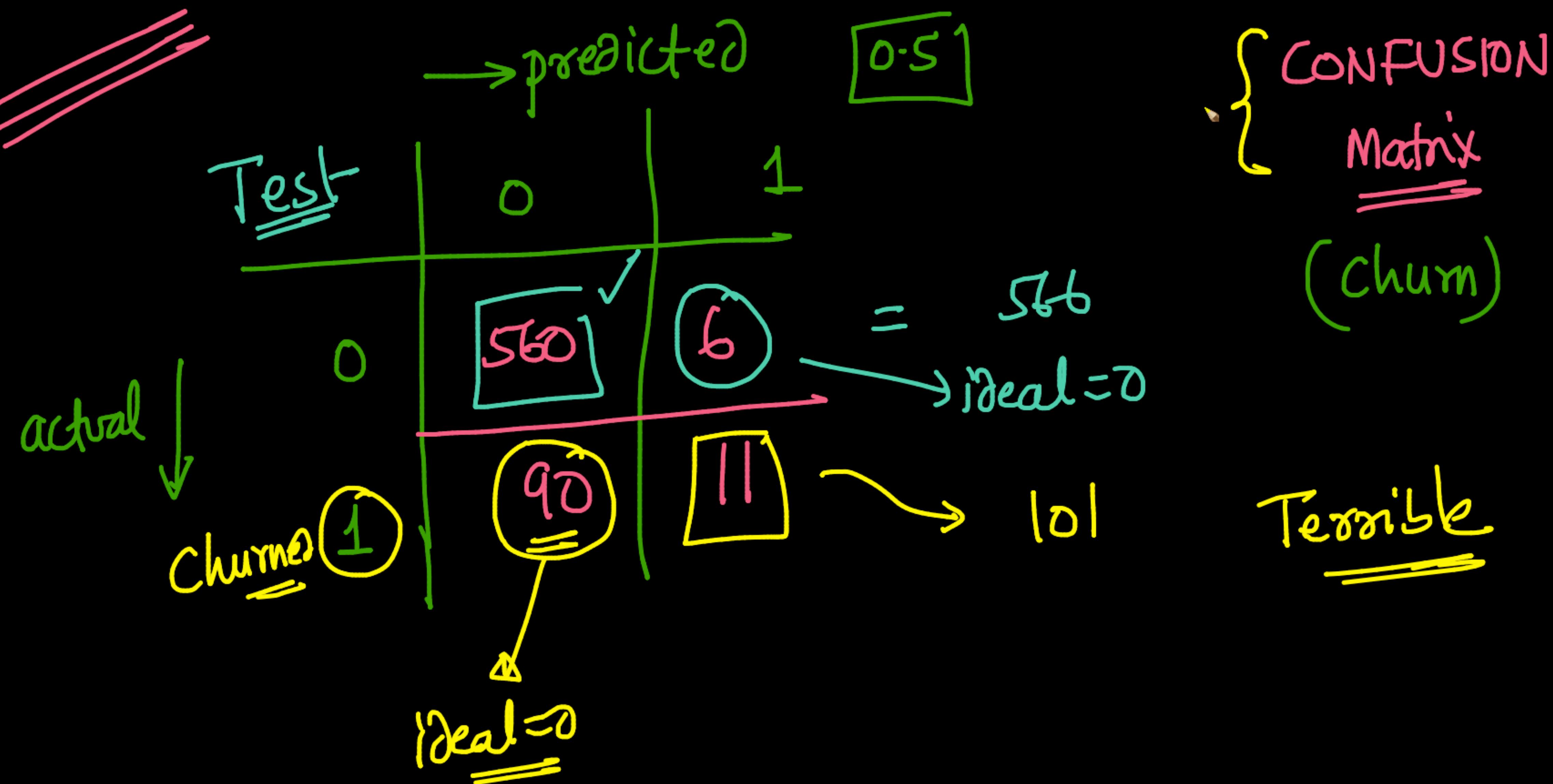
[] np.argmax(val_scores)

[] val_scores[0]

0.856071964017991

Reconnect

17 / 17



+ Code + Text

```
[ ] from sklearn.metrics import accuracy_score  
  
{x} print(f"Accuracy : {accuracy_score(y_test, y_pred)*100}%")  
  
Accuracy : 85.6071964017991%
```

```
[ ] from sklearn.metrics import confusion_matrix  
  
confusion_matrix(y_test, y_pred)  
  
array([[560,   6],  
       [ 90,  11]])
```

```
[ ] # many class-1 points classified as class-0 => useless model
```

```
[ ] # how to fix the model?  
  
# Hyper-param tuning  
from sklearn.pipeline import make_pipeline  
train_scores = []  
val_scores = []  
scaler = StandardScaler()
```

+ Code + Text

Reconnect



LogisticRegression(C=0.0142836737608913)))

[] test_score = scaled_lr.score(X_test, y_test)

print(test score)

{x} y_pred = scaled_lr.predict(X_test)

0.856071964017991

[] from sklearn.metrics import accuracy_score

print(f"Accuracy : {accuracy_score(y_test, y_pred)*100}%")

Accuracy : 85.6071964017991%

<> [] from sklearn.metrics import confusion_matrix

confusion_matrix(y_test, y_pred)

<> array([[560, 6],
[90, 11]])

>- [] # many class-1 points classified as class-0 => useless model

+ Code + Text

Reconnect



```
[ ] from sklearn.metrics import accuracy_score  
print(f"Accuracy : {accuracy_score(y_test, y_pred)*100}%")
```

Accuracy : 85.6071964017991%

y_i \hat{y}_i

```
[ ] from sklearn.metrics import confusion_matrix  
  
confusion_matrix(y_test, y_pred)
```

```
array([[560,    6],  
       [ 90,   11]])
```

```
[ ] # many class-1 points classified as class-0 => useless model
```

```
[ ] # how to fix the model?
```

```
# Hyper-param tuning  
from sklearn.pipeline import make_pipeline  
train_scores = []  
val_scores = []  
scaler = StandardScaler()
```

+ Code + Text

Reconnect



```
[ ] from sklearn.metrics import accuracy_score  
  
print(f"Accuracy : {accuracy_score(y_test, y_pred)*100}%")
```

Accuracy : 85.6071964017991%

```
[ ] from sklearn.metrics import confusion_matrix
```

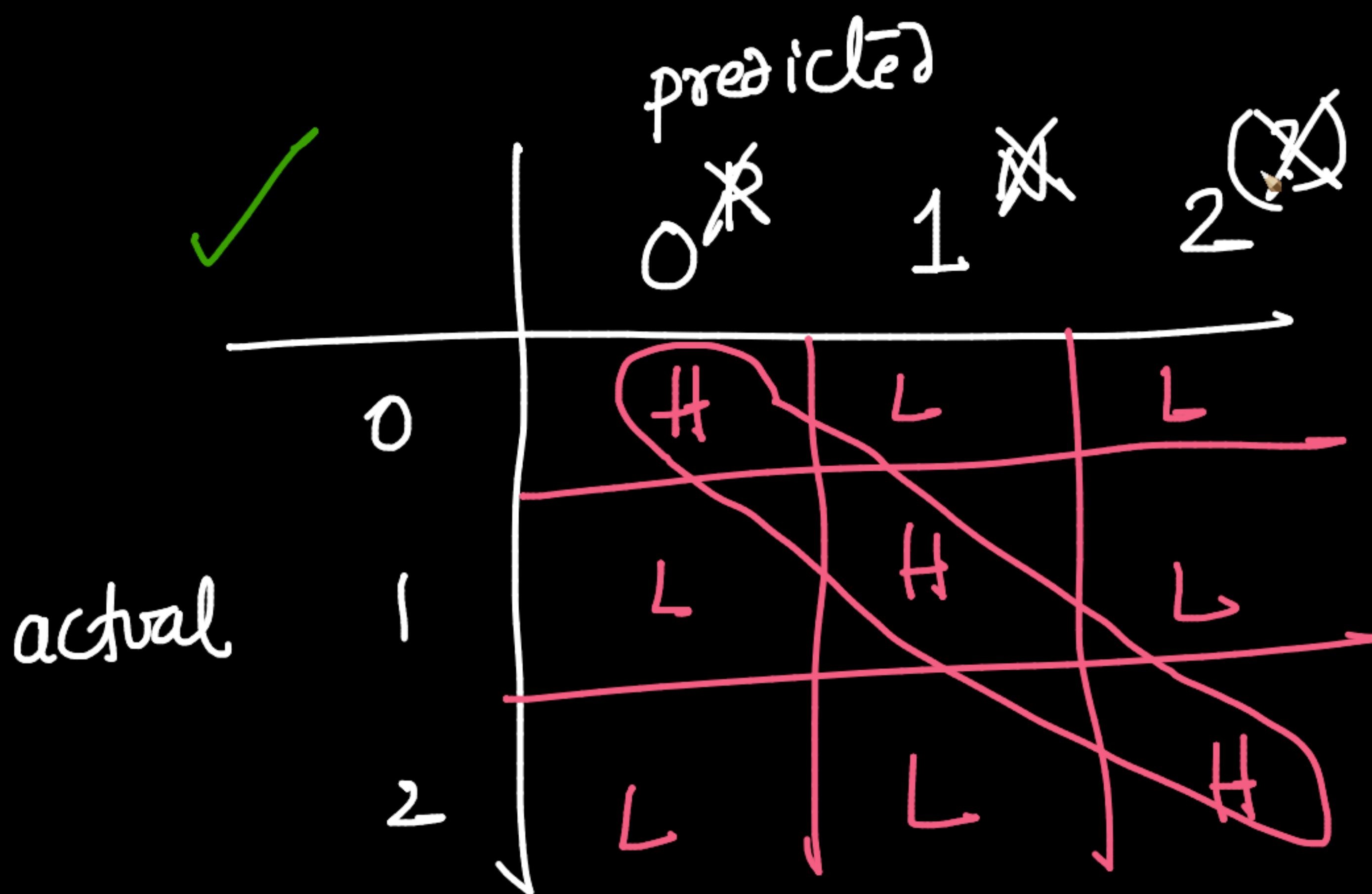
```
confusion_matrix(y_test, y_pred)
```

```
array([[560,  6],  
       [ 90, 11]])
```

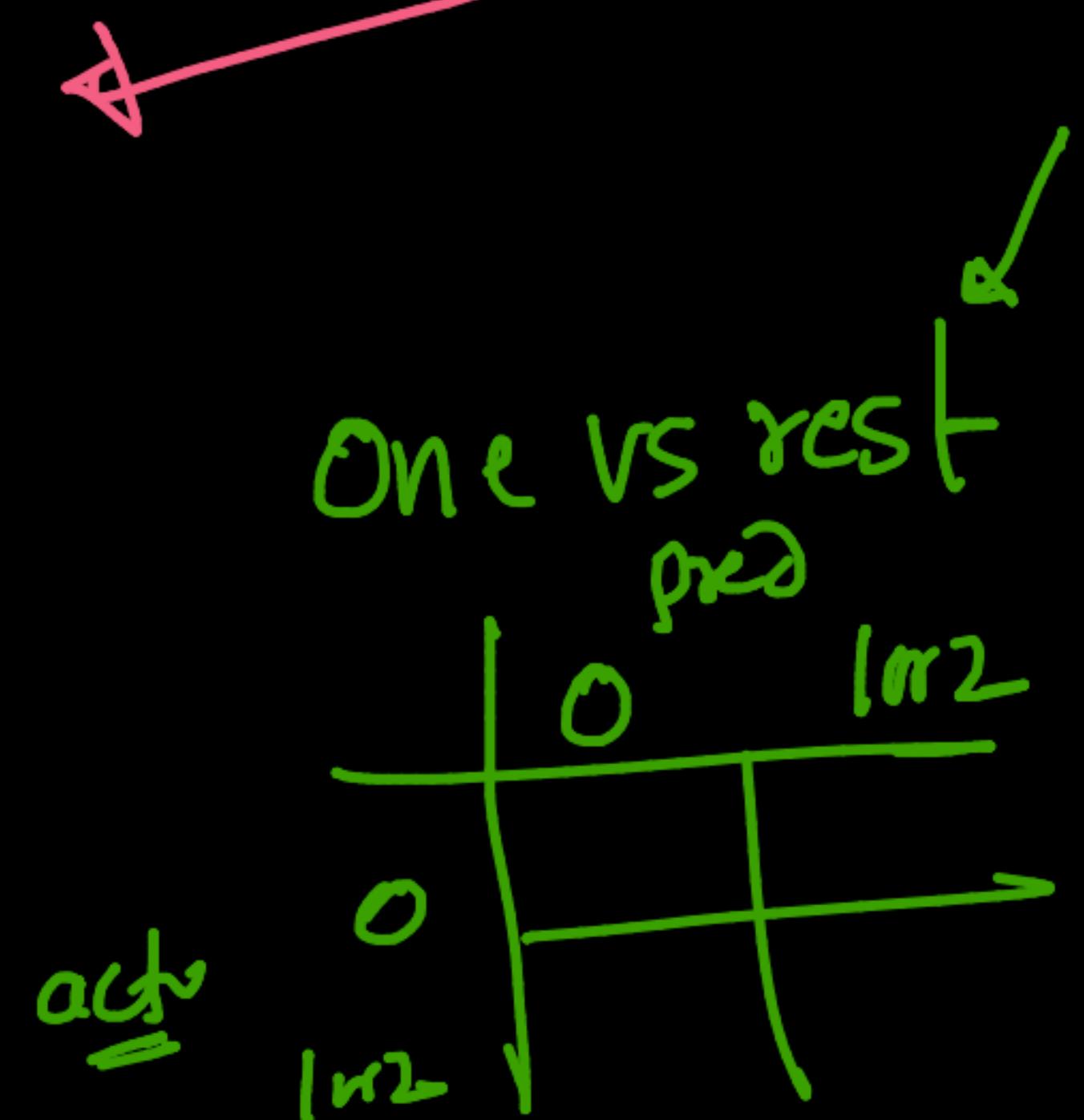
```
[ ] # many class-1 points classified as class-0 => useless model
```

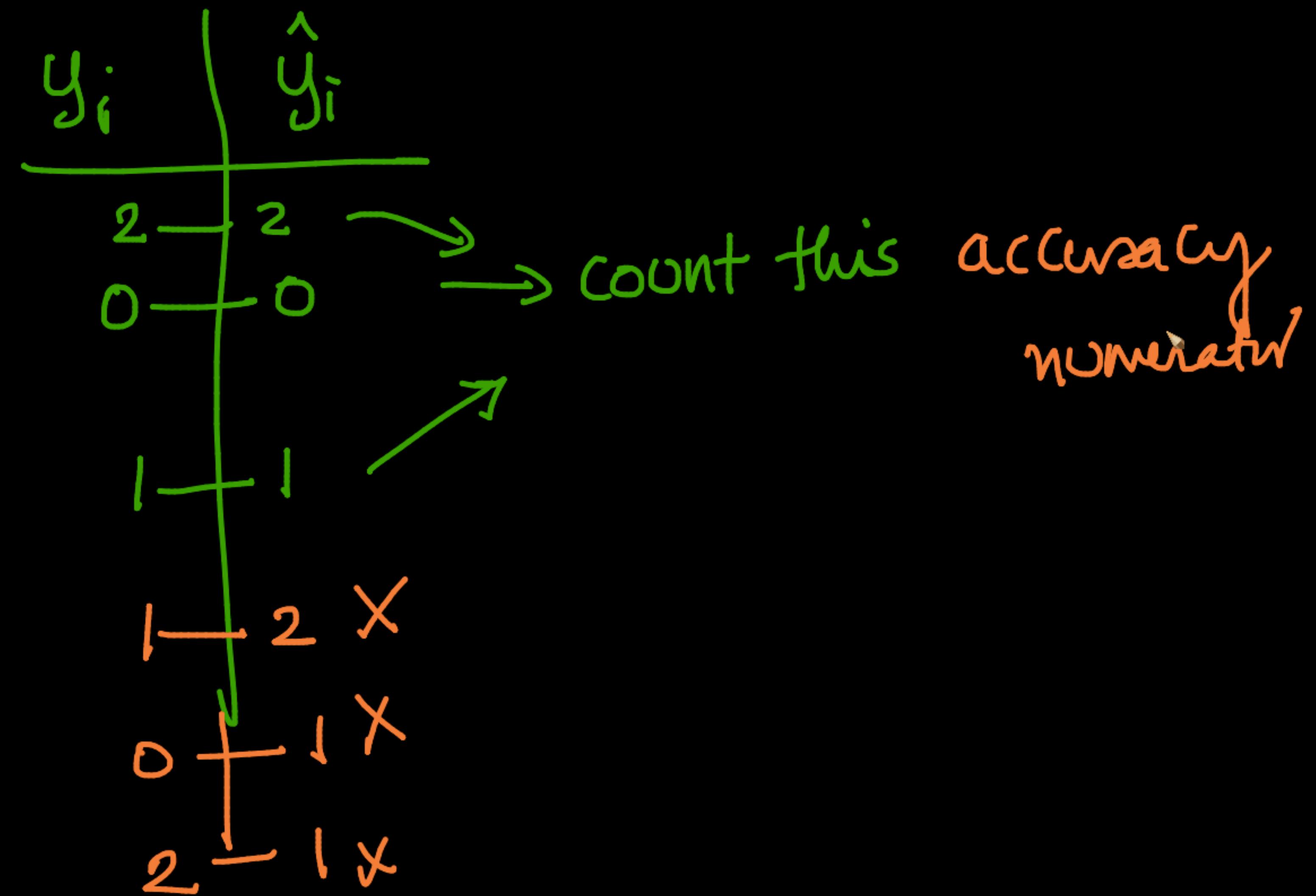
```
[ ] # how to fix the model?
```

```
# Hyper-param tuning  
from sklearn.pipeline import make_pipeline  
train_scores = []  
val_scores = []  
scaler = StandardScaler()
```



CONF - Matrix
3 classes





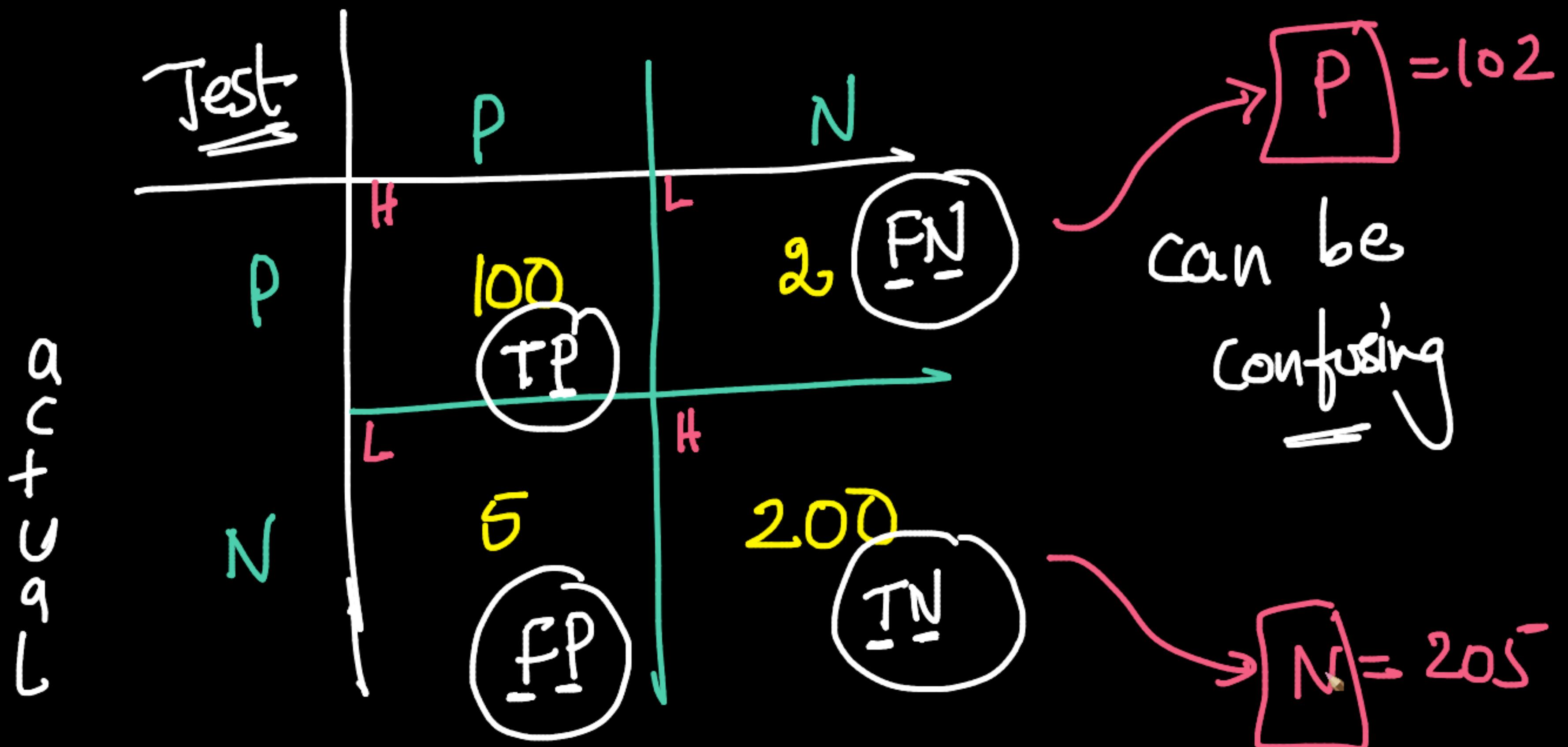
Confusion-Matrix in Sklearn

↳ one vs rest 
=

CONFUSION - Matrix

predicted (Model)

~~binary
classifi-~~



= Case-study



fix



+ Code + Text

Reconnect



```
[ ] print(f"Accuracy : {accuracy_score(y_test, y_pred)*100}%")
```

Accuracy : 85.6071964017991%

```
[ ] from sklearn.metrics import confusion_matrix
```

```
confusion_matrix(y_test, y_pred)
```

```
array([[560,    6],  
       [ 90,   11]]) ] -v. bad
```

```
[ ] # many class-1 points classified as class-0 => useless model
```

```
[ ] # how to fix the model?
```

```
# Hyper-param tuning  
from sklearn.pipeline import make_pipeline  
train_scores = []  
val_scores = []  
scaler = StandardScaler()  
  
for la in np.arange(0.01, 100.0, 5):  
    scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la, class_weight={0:0.1, 1:0.9}))
```

pred	0	1
act	low	high
0	high	low
1	low	high



+ Code + Text

```
from sklearn.metrics import accuracy_score  
[ ]  
print(f"Accuracy : {accuracy_score(y_test, y_pred)*100}%")
```

{x}
Accuracy : 85.6071964017991%

```
[ ] from sklearn.metrics import confusion_matrix  
  
confusion_matrix(y_test, y_pred)  
  
array([[560,    6,  
       1,  90],  
      [ 90, 11]])
```

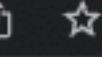
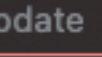
```
[ ] # many class-1 points classified as class-0 => useless model
```

```
[ ] # how to fix the model?  
  
# Hyper-param tuning  
from sklearn.pipeline import make_pipeline  
train_scores = []  
val_scores = []  
scaler = StandardScaler()  
  
for la in np.arange(0.01, 100.0, 5):
```

Logistic Regression.ipynb - Colab | sklearn.linear_model.LogisticR... | Classifier comparison – scikit-learn | Confusion matrix - Wikipedia | +

colab.research.google.com/drive/1G_D39wa40kVahULehDS4lWhB_PXhP-xJ#scrollTo=X2vZAFjWc-8E

+ Code + Text

Reconnect |   

```
from sklearn.metrics import accuracy_score
[ ]
print(f"Accuracy : {accuracy_score(y_test, y_pred)*100}%")
```

{x}

```
Accuracy : 85.6071964017991%
```

[] from sklearn.metrics import confusion_matrix

```
confusion_matrix(y_test, y_pred)
```

```
array([[560,   6],
       [ 90,  11]])
```

[] # many class-1 points classified as class-0 => useless model

[] # how to fix the model?

```
# Hyper-param tuning
from sklearn.pipeline import make_pipeline
train_scores = []
val_scores = []
scaler = StandardScaler()

for la in np.arange(0.01, 100.0, 5):
```

fixes → imbalance correction
↓
weights

Logistic Regression.ipynb - Colab | sklearn.linear_model.LogisticR... | Classifier comparison – scikit-learn.org | Confusion matrix - Wikipedia

scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html

Install User Guide API Examples Community More ▾

class_weight 2/3 Go

Prev Up Next

scikit-learn 1.0.2

Other versions

Please cite us if you use the software.

sklearn.linear_model.LogisticRegression

Examples using

sklearn.linear_model.LogisticRe

{ 0 : ↓
1 : ↑

sklearn.linear_model.LogisticRegression

```
class sklearn.linear_model.LogisticRegression(penalty='l2', *, dual=False, tol=0.0001, C=1.0, fit_intercept=True,  
intercept_scaling=1, class_weight=None, random_state=None, solver='lbfgs', max_iter=100, multi_class='auto', verbose=0,  
warm_start=False, n_jobs=None, l1_ratio=None)
```

[source]

Logistic Regression (aka logit, MaxEnt) classifier.

In the multiclass case, the training algorithm uses the one-vs-rest (OvR) scheme if the 'multi_class' option is set to 'ovr', and uses the cross-entropy loss if the 'multi_class' option is set to 'multinomial'. (Currently the 'multinomial' option is supported only by the 'lbfgs', 'sag', 'saga' and 'newton-cg' solvers.)

This class implements regularized logistic regression using the 'liblinear' library, 'newton-cg', 'sag', 'saga' and 'lbfgs' solvers. **Note that regularization is applied by default.** It can handle both dense and sparse input. Use C-ordered arrays or CSR matrices containing 64-bit floats for optimal performance; any other input format will be converted (and copied).

The 'newton-cg', 'sag', and 'lbfgs' solvers support only L2 regularization with primal formulation, or no regularization. The 'liblinear' solver supports both L1 and L2 regularization, with a dual formulation only for the L2 penalty. The Elastic-Net regularization is only supported by the 'saga' solver.

Read more in the [User Guide](#).

Parameters: `penalty : {'l1', 'l2', 'elasticnet', 'none'}, default='l2'`

Specify the norm of the penalty:

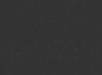
- 'none' : no penalty is added;
- 'l2' : add a L2 penalty term and it is the default choice;
- 'l1' : add a L1 penalty term;
- 'elasticnet' : both L1 and L2 penalty terms are added.



31 / 31

Toggle Menu

. See the parameter solver below, to know

 + Code + Text

$$\begin{aligned} \text{c-0} &: 0.1 \\ \text{c-1} &: 0.9 \end{aligned}$$

```
[ ] # how to fix the model?

# Hyper-param tuning
from sklearn.pipeline import make_pipeline
train_scores = []
val_scores = []
scaler = StandardScaler()

for la in np.arange(0.01, 100.0, 5):
    scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la, class_weight={0:0.1, 1:0.9}))
    scaled_lr.fit(X_train, y_train)
    train_score = scaled_lr.score(X_train, y_train)
    val_score = scaled_lr.score(X_val, y_val)
    train_scores.append(train_score)
    val_scores.append(val_score)

plt.figure()
plt.plot(list(np.arange(0.01, 100.0, 5)), train_scores, label="train")
plt.plot(list(np.arange(0.01, 100.0, 5)), val_scores, label="val")
plt.legend(loc='lower right')
plt.xlabel("lambda")
plt.ylabel("Score")
plt.grid()
plt.show()
```

$$\sum_{i=1}^n \log\text{-loss}_i(w_i)$$

+ Code + Text

Reconnect



0.856071964017991

```
[ ] from sklearn.metrics import accuracy_score  
{x}  
print(f"Accuracy : {accuracy_score(y_test, y_pred)*100}%")
```

Accuracy : 85.6071964017991%



```
▶ from sklearn.metrics import confusion_matrix  
  
confusion_matrix(y_test, y_pred)  
  
array([[560, 6],  
       [ 90, 11]])  
  
[ ] # many class-1 points classified as class-0 => useless model
```

how to fix the model?

```
# Hyper-param tuning  
from sklearn.pipeline import make_pipeline  
train_scores = []  
val_scores = []
```

[] # many class-1 points classified as class-0 => useless model

[] # how to fix the model?

Hyper-param tuning

```
from sklearn.pipeline import make_pipeline
train_scores = []
val_scores = []
scaler = StandardScaler()

for la in np.arange(0.01, 100.0, 5):
    scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la, class_weight={0:0.1, 1:0.9}))
    scaled_lr.fit(X_train, y_train)
    train_score = scaled_lr.score(X_train, y_train)
    val_score = scaled_lr.score(X_val, y_val)
    train_scores.append(train_score)
    val_scores.append(val_score)

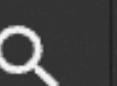
plt.figure()
plt.plot(list(np.arange(0.01, 100.0, 5)), train_scores, label="train")
plt.plot(list(np.arange(0.01, 100.0, 5)), val_scores, label="val")
plt.legend(loc='lower right')
plt.xlabel("lambda")
plt.ylabel("Score")
plt.grid()
```

bias-variance

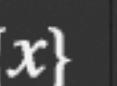
hyper-param



+ Code



[] # many class-1 points classified as class-0 => useless model



[] # how to fix the model?



```
# Hyper-param tuning
from sklearn.pipeline import make_pipeline
train_scores = []
val_scores = []
scaler = StandardScaler()

for la in np.arange(0.01, 100.0, 5):
    scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la, class_weight={0:0.1, 1:0.9}))
    scaled_lr.fit(X_train, y_train)
    train_score = scaled_lr.score(X_train, y_train)
    val_score = scaled_lr.score(X_val, y_val)
    train_scores.append(train_score)
    val_scores.append(val_score)

plt.figure()
plt.plot(list(np.arange(0.01, 100.0, 5)), train_scores, label="train")
plt.plot(list(np.arange(0.01, 100.0, 5)), val_scores, label="val")
plt.legend(loc='lower right')
plt.xlabel("lambda")
plt.ylabel("Score")
plt.grid()
```

15:85

+ Code + Text

Reconnect



```
[ ] # many class-1 points classified as class-0 => useless model

[ ] # how to fix the model?

{x}
# Hyper-pram tuning
from sklearn.pipeline import make_pipeline
train_scores = []
val_scores = []
scaler = StandardScaler()

for la in np.arange(0.01, 100.0, 5):
    scaled_lr = make_pipeline( scaler, LogisticRegression(C=1/la , class_weight={ 0:0.1, 1:0.9 })) 
    scaled_lr.fit(X_train, y_train)
    train_score = scaled_lr.score(X_train, y_train)
    val_score = scaled_lr.score(X_val, y_val)
    train_scores.append(train_score)
    val_scores.append(val_score)

plt.figure()
plt.plot(list(np.arange(0.01, 100.0, 5)), train_scores, label="train")
plt.plot(list(np.arange(0.01, 100.0, 5)), val_scores, label="val")
plt.legend(loc='lower right')
plt.xlabel("lambda")
plt.ylabel("Score")
plt.grid()
```

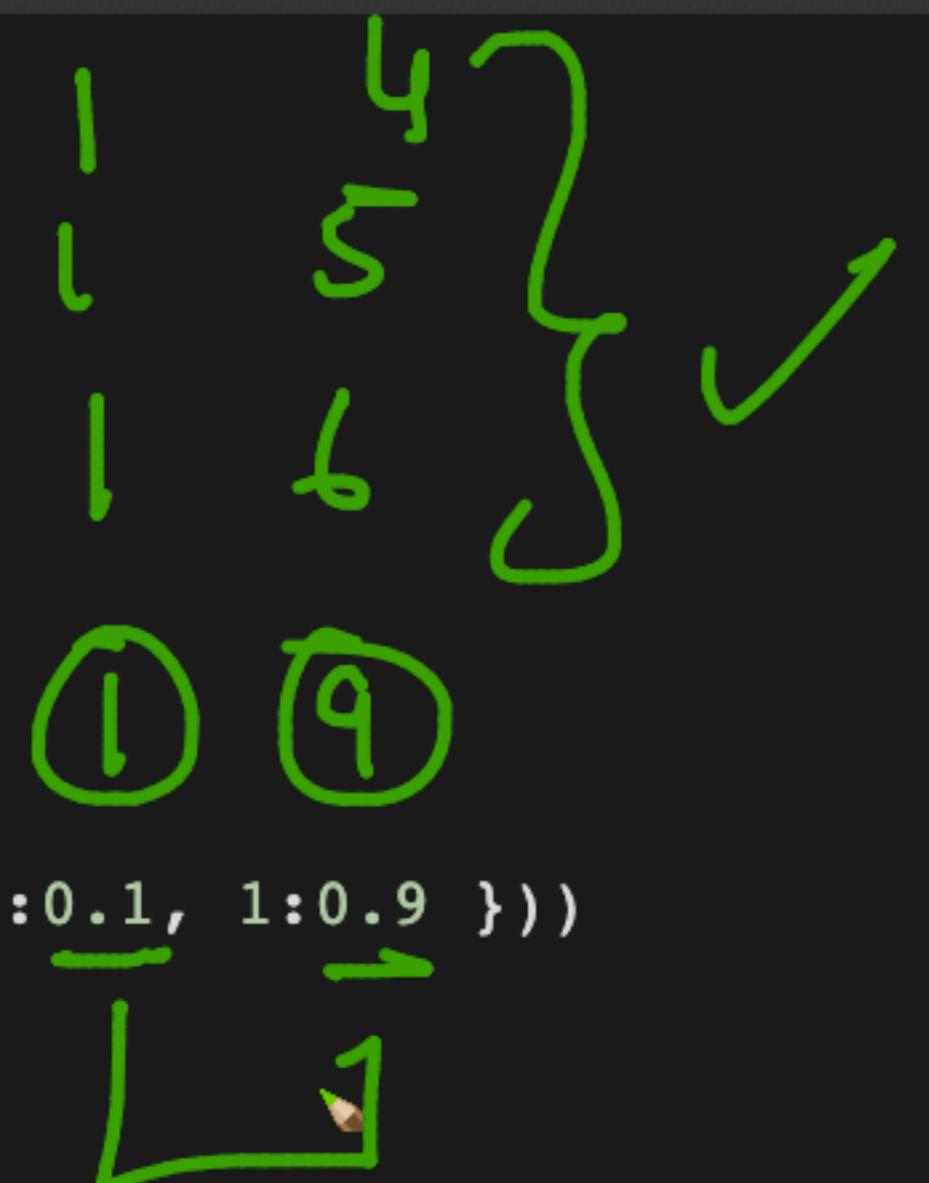
[] # many class-1 points classified as class-0 => useless model

[] # how to fix the model?

```
# Hyper-param tuning
from sklearn.pipeline import make_pipeline
train_scores = []
val_scores = []
scaler = StandardScaler()

for la in np.arange(0.01, 100.0, 5):
    scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la, class_weight={0:0.1, 1:0.9}))
    scaled_lr.fit(X_train, y_train)
    train_score = scaled_lr.score(X_train, y_train)
    val_score = scaled_lr.score(X_val, y_val)
    train_scores.append(train_score)
    val_scores.append(val_score)

plt.figure()
plt.plot(list(np.arange(0.01, 100.0, 5)), train_scores, label="train")
plt.plot(list(np.arange(0.01, 100.0, 5)), val_scores, label="val")
plt.legend(loc='lower right')
plt.xlabel("lambda")
plt.ylabel("Score")
plt.grid()
```



+ Code + Text

Reconnect  Update

```
[ ] # many class-1 points classified as class-0 => useless model

[ ] # how to fix the model?

{x}
# Hyper-param tuning
from sklearn.pipeline import make_pipeline
train_scores = []
val_scores = []
scaler = StandardScaler()

for la in np.arange(0.01, 100.0, 5):
    scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la, class_weight={0:0.1, 1:0.9}))
    scaled_lr.fit(X_train, y_train)
    train_score = scaled_lr.score(X_train, y_train)
    val_score = scaled_lr.score(X_val, y_val)
    train_scores.append(train_score)
    val_scores.append(val_score)

plt.figure()
plt.plot(list(np.arange(0.01, 100.0, 5)), train_scores, label="train")
plt.plot(list(np.arange(0.01, 100.0, 5)), val_scores, label="val")
plt.legend(loc='lower right')
plt.xlabel("lambda")
plt.ylabel("Score")
plt.grid()
```

ex: tune ratios of weights

↪ $\frac{1}{9}$

Logistic Regression.ipynb - Colab | sklearn.linear_model.LogisticR... | Classifier comparison – scikit-learn | Confusion matrix - Wikipedia | +

colab.research.google.com/drive/1G_D39wa40kVahULehDS4lWhB_PXhP-xJ#scrollTo=xaQEhH9_uOvt

+ Code + Text

RAM Disk

1s

{x}

□

→ log-loss

acc

Overfit Underfit

```
plt.figure()
plt.plot(list(np.arange(0.01, 100.0, 5)), train_scores, label="train")
plt.plot(list(np.arange(0.01, 100.0, 5)), val_scores, label="val")
plt.legend(loc='lower right')
plt.xlabel("lambda")
plt.ylabel("Score")
plt.grid()
plt.show()
```

39 / 39

+ Code + Text

```
[119] print(f"Accuracy : {accuracy score(y test, y pred)*100}%")
```

Accuracy : 70.4647676161919%

```
from sklearn.metrics import confusion_matrix  
  
confusion_matrix(y_test, y_pred)  
# Better for class-1, but worse for class-0  
# Any ideas on how to solve?  
  
array([[383, 183],  
       [ 14,  87]])
```

[]

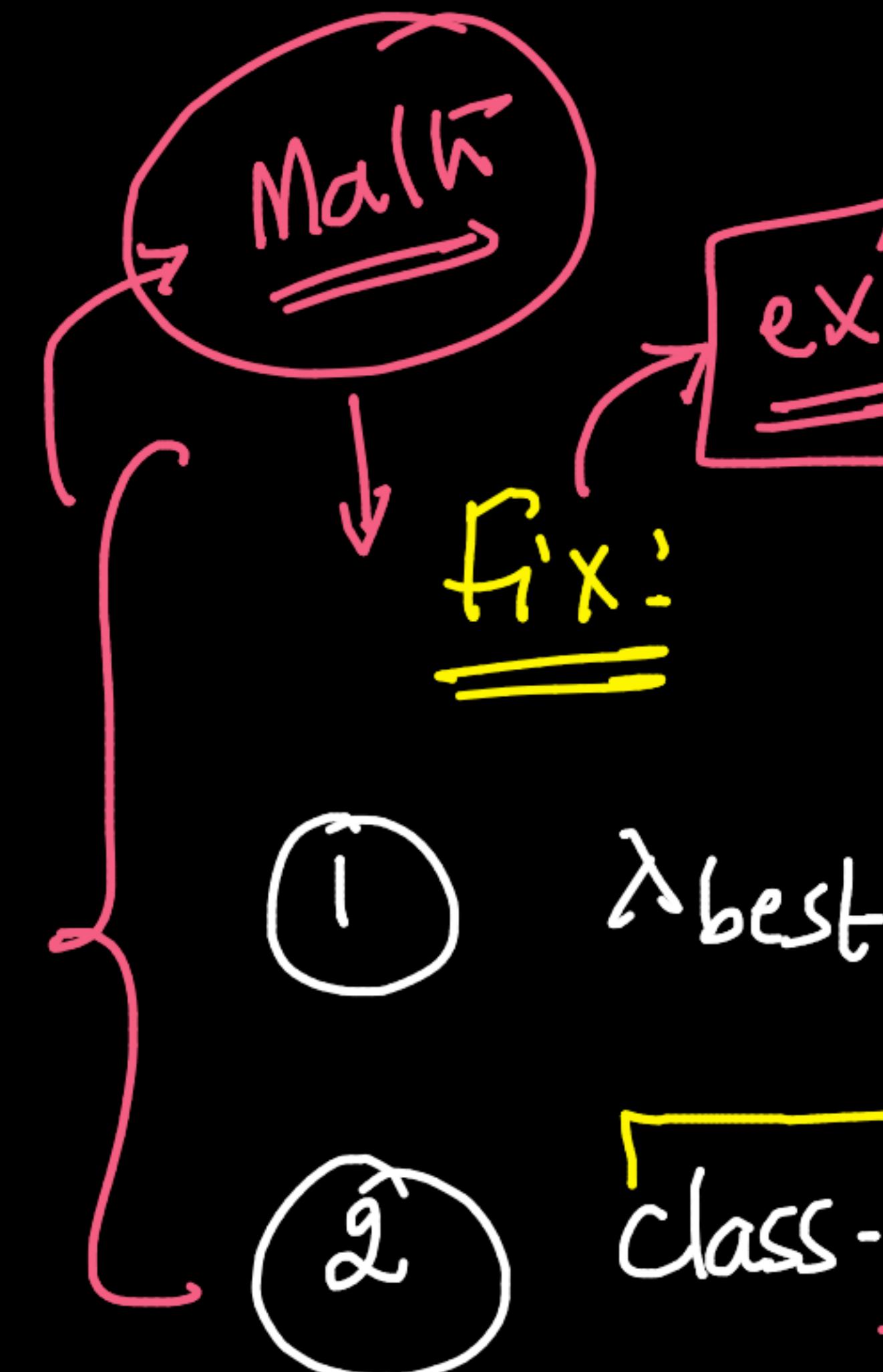
pred		
0	1	
actual	0	14
0	383	183
1	87	✓

d0 d1

1 9

$$\frac{1}{2}; \frac{1}{5}; \frac{1}{6}; \frac{1}{9}; \frac{1}{10}; \dots$$

0	0	1
0	H	L
1	L	H

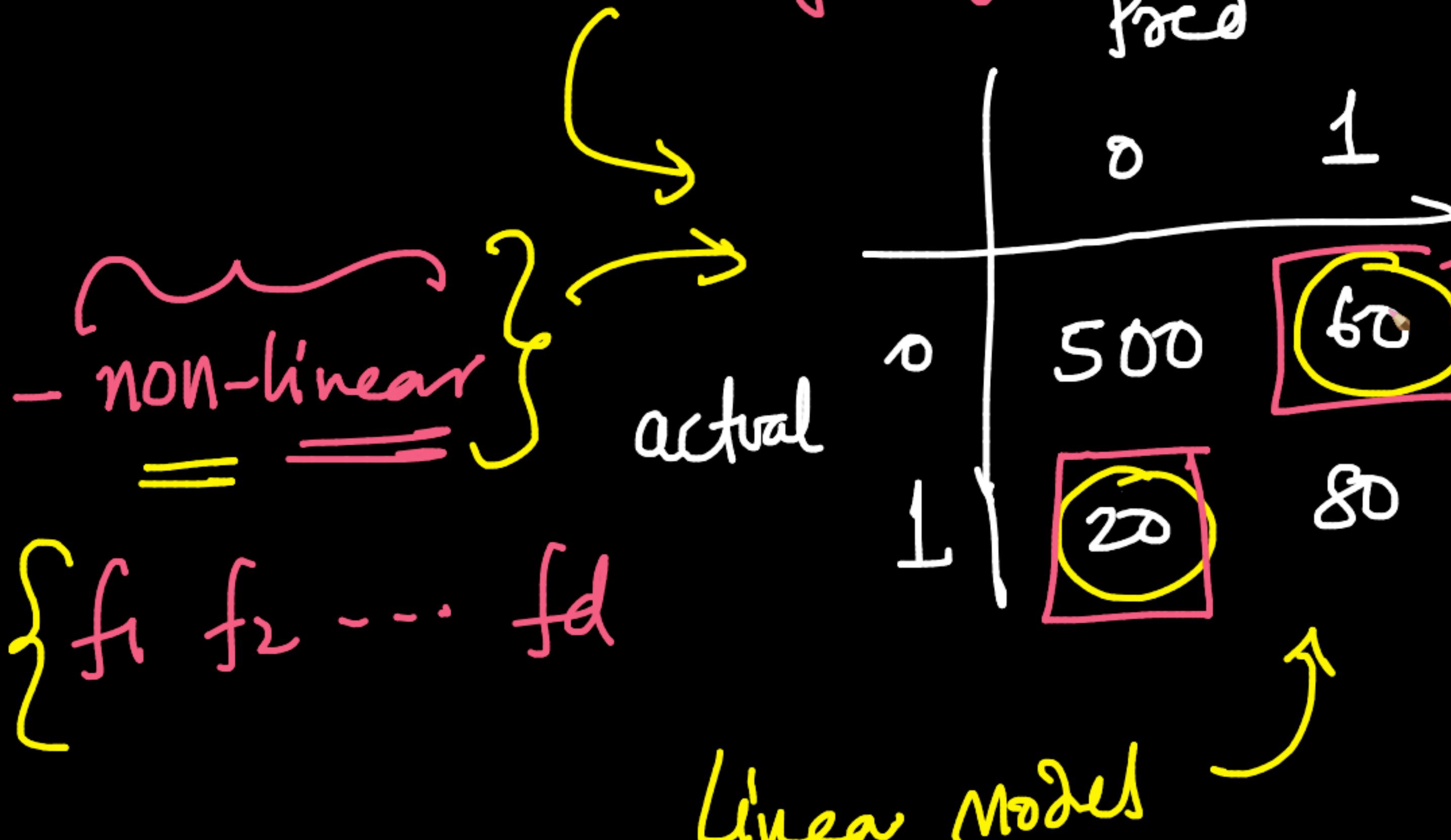


class-weights ratio $\frac{1}{a}$ ✓
hyperparameter

log-loss λ_{best}

$$\frac{1}{a}$$

Not-separable linearly
fix - everything



How can you
further improve

- - outliers
 - look @ FP & FN
- 0.5 to 0.4 ...
hyper-param

no rule of thumb

1:5
0.1 - 0.9

tweak
(hyperparam)

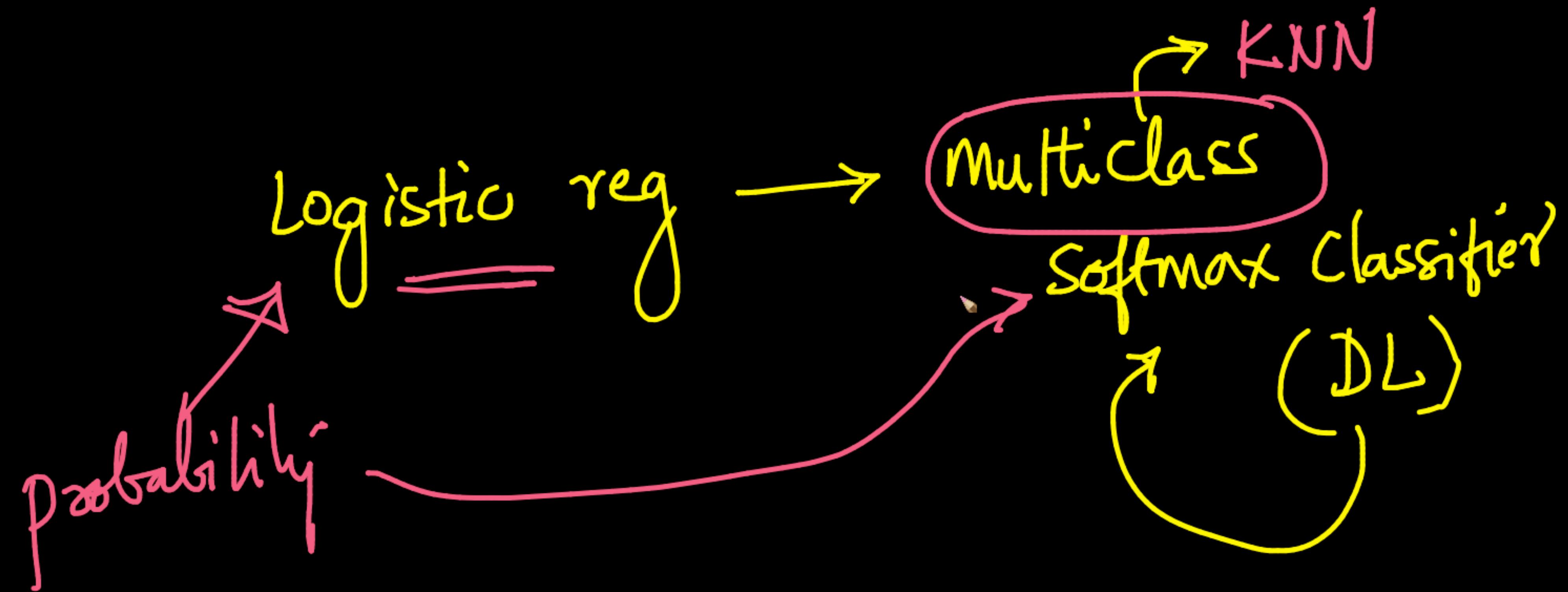
← fine-tune

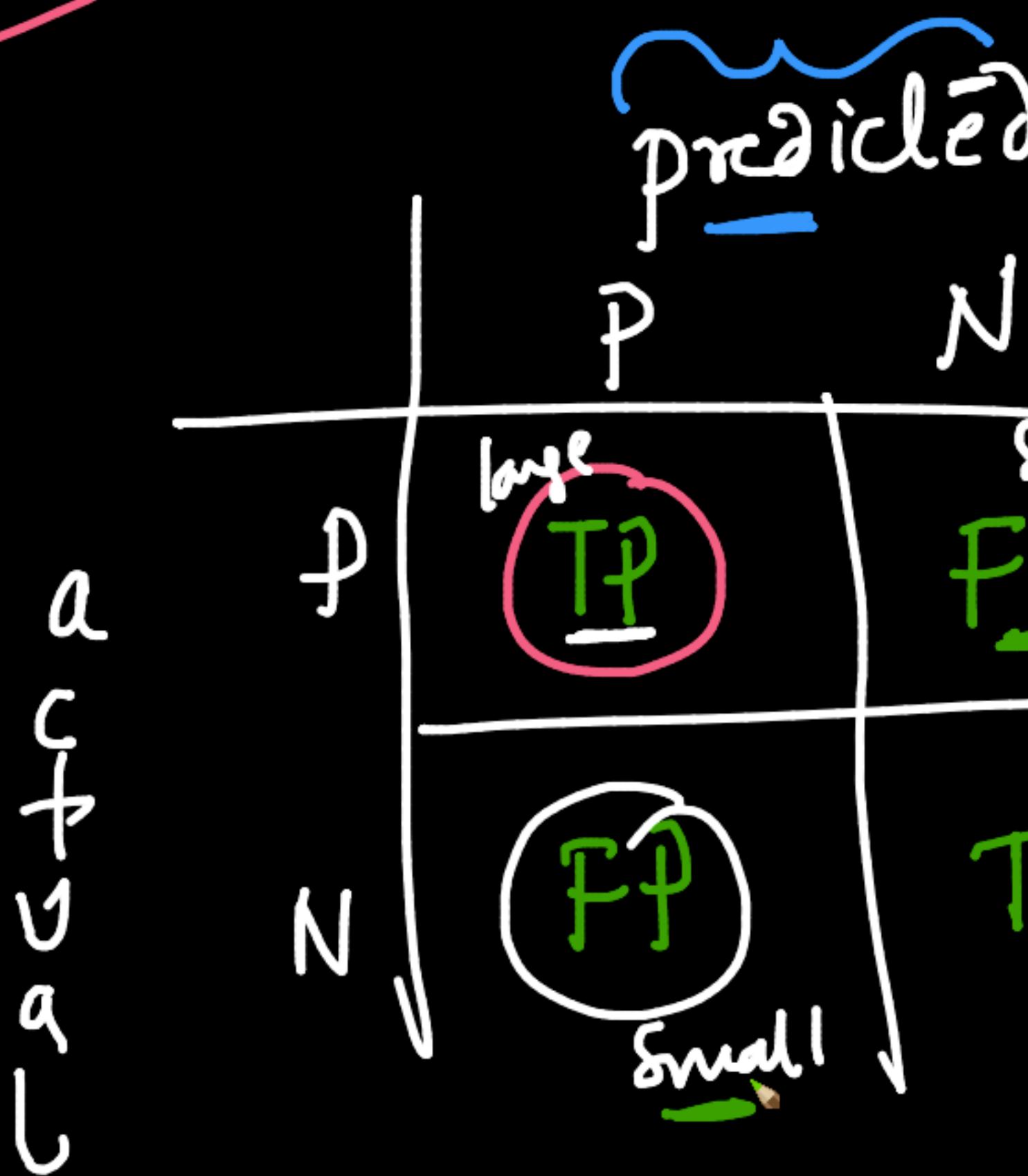
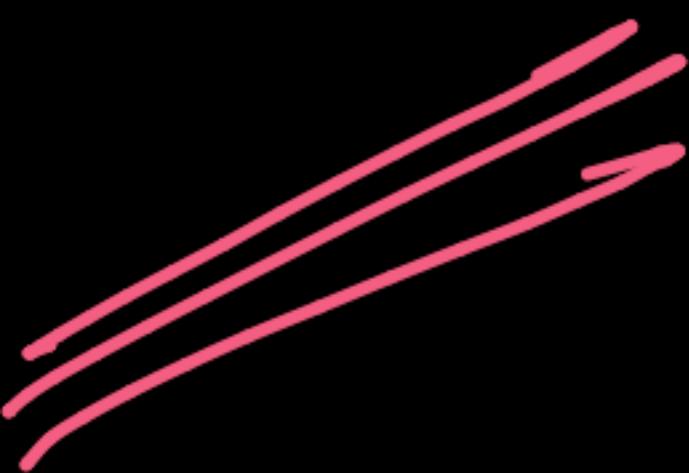
DL do

15 : 85

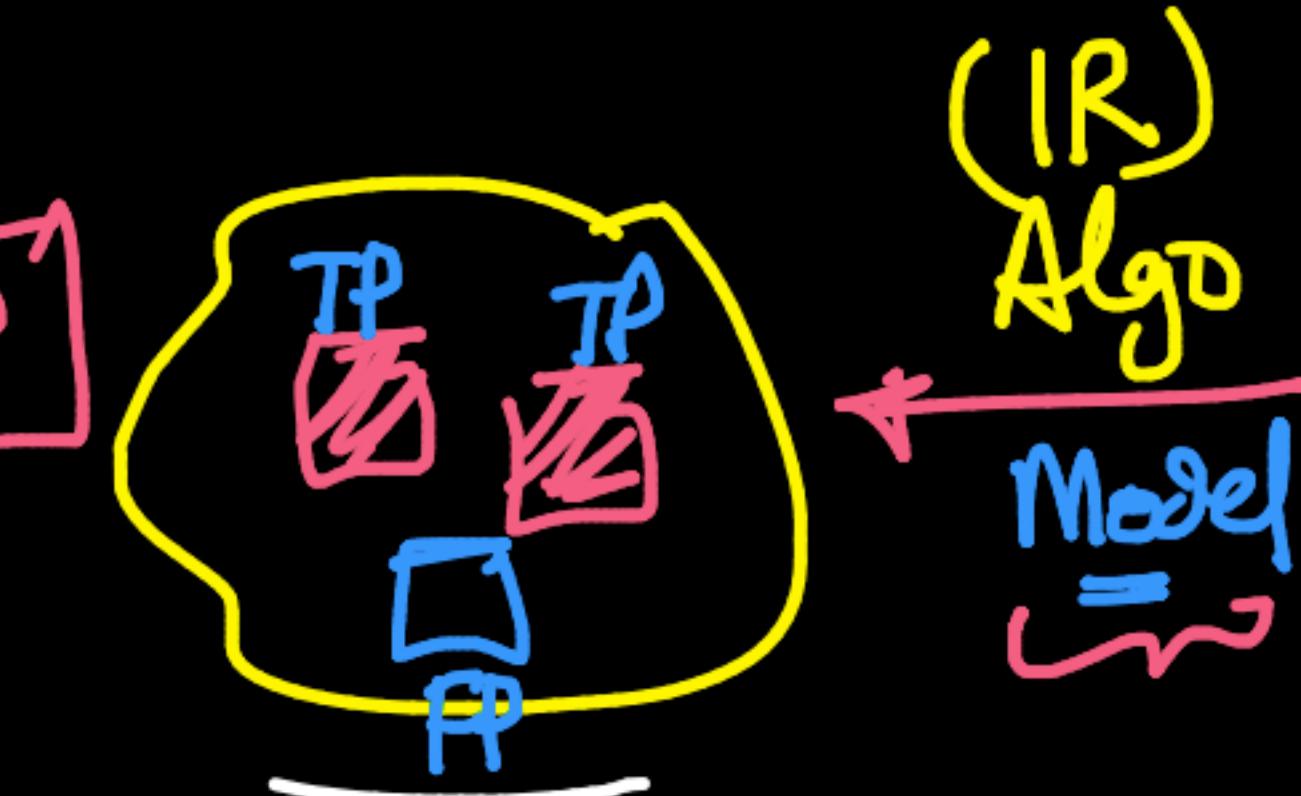
1 : 6

Wi : 6 : 1





→ { Information - retrieval:



✓ Precision:

$$\frac{2}{3} = \frac{TP}{TP+FP}$$

Recall:

$$\frac{2}{4} = \frac{TP}{TP+FN}$$

[Permanent link](#)[Page information](#)[Cite this page](#)[Wikidata item](#)[Print/export](#)[Download as PDF](#)[Printable version](#)

Languages



العربية

Deutsch

Español

Français

한국어

Italiano

日本語

Português

中文

文 A 5 more

Edit links

5 References

Example [edit]

Given a sample of 12 individuals, 8 that have been diagnosed with cancer and 4 that are cancer-free, where individuals with cancer belong to class 1 (positive) and non-cancer individuals belong to class 0 (negative), we can display that data as follows:

Individual Number	1	2	3	4	5	6	7	8	9	10	11	12
Actual Classification	1	1	1	1	1	1	1	1	0	0	0	0

Assume that we have a classifier that distinguishes between individuals with and without cancer in some way, we can take the 12 individuals and run them through the classifier. The classifier then makes 9 accurate predictions and misses 3: 2 individuals with cancer wrongly predicted as being cancer-free (sample 1 and 2), and 1 person without cancer that is wrongly predicted to have cancer (sample 9).

Individual Number	1	2	3	4	5	6	7	8	9	10	11	12
Actual Classification	1	1	1	1	1	1	1	1	0	0	0	0
Predicted Classification	0	0	1	1	1	1	1	1	1	0	0	0

Notice, that if we compare the actual classification set to the predicted classification set, there are 4 different outcomes that could result in any particular column. One, if the actual classification is positive and the predicted classification is positive (1,1), this is called a true positive result because the positive sample was correctly identified by the classifier. Two, if the actual classification is positive and the predicted classification is negative (1,0), this is called a false negative result because the positive sample is incorrectly identified by the classifier as being negative. This classification is positive (0,1), this is called a false positive result because the negative

A test result which wrongly indicates that a particular condition or attribute is absent

sensitivity, recall hit rate, or true positive rate (TPR)

$$\text{TPR} = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 - \text{FNR}$$

specificity, selectivity or true negative rate (TNR)

$$\text{TNR} = \frac{\text{TN}}{\text{N}} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$$

precision or positive predictive value (PPV)

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} = 1 - \text{FDR}$$

negative predictive value (NPV)

$$\text{NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}} = 1 - \text{FOR}$$

miss rate or false negative rate (FNR)

$$\text{FNR} = \frac{\text{FN}}{\text{P}} = \frac{\text{FN}}{\text{FN} + \text{TP}} = 1 - \text{TPR}$$

false-out or false positive rate (FPR)

$$\text{FPR} = \frac{\text{FP}}{\text{N}} = \frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{TNR}$$

false discovery rate (FDR)

$$\text{FDR} = \frac{\text{FP}}{\text{FP} + \text{TP}} = 1 - \text{PPV}$$

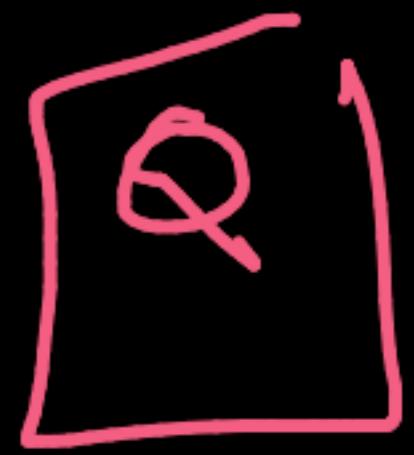
false omission rate (FOR)

$$\text{FOR} = \frac{\text{FN}}{\text{FN} + \text{TN}} = 1 - \text{NPV}$$

Positive likelihood ratio (LR+)

$$\text{LR+} = \frac{\text{TPR}}{\text{FPR}}$$

Negative likelihood ratio (LR-)



Precision : ↑
↓ Recall :

Tradeoff:

Pr & Recall

e.g.: covid diagnosis
RT-PCR: Recall {
v-high }
99% → high Precision

Precision : $\frac{TP}{(TP+FP)}$

Recall : $\frac{TP}{(TP+FN)}$

→ one number
metric

Pr↑ & re↑

$$\left\{ \begin{array}{l} 0-1 \\ 0-1 \end{array} \right.$$

$$\text{avg} = \frac{\uparrow p_{\gamma} + \text{Re} \uparrow}{z} \quad - \text{arithmetic Mean} \uparrow \alpha p_{\gamma}$$

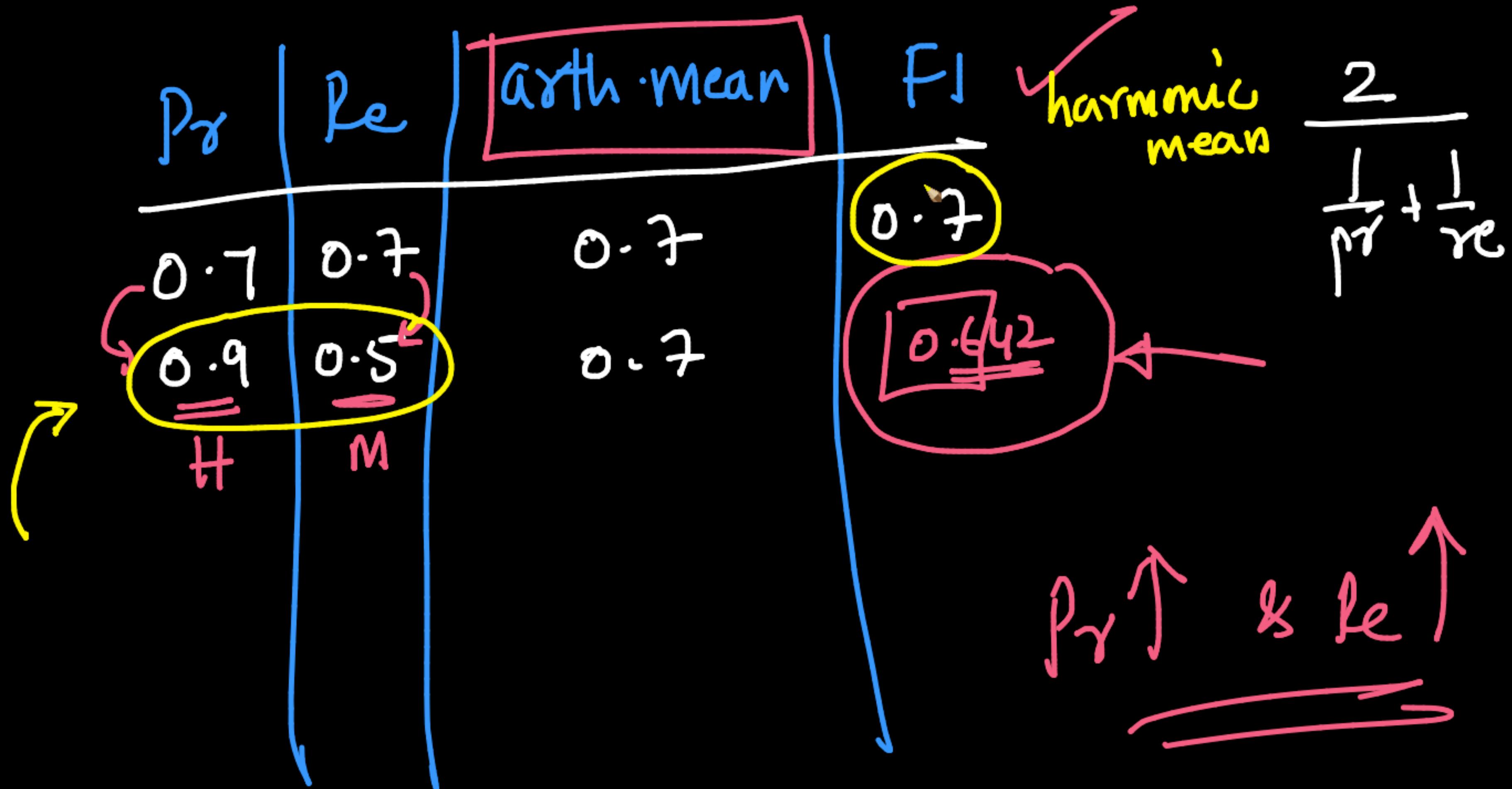
$\alpha \text{ Re}$

↑↑

$$\left\{ \begin{array}{l} \overline{f_1} \\ \overline{s_{\text{ave}}} \end{array} \right.$$
$$\frac{2}{\left(\frac{1}{\uparrow p_{\gamma}} + \frac{1}{\text{Re} \uparrow} \right)} : \text{harmonic mean}$$

=====

$$\left\{ \begin{array}{l} \text{Re} \propto f_1 \\ p_{\gamma} \propto f_1 \end{array} \right.$$



$R_Y \uparrow$ & $R_E \uparrow$

$$\left\{ \begin{array}{l} f_1 \\ \hline \end{array} \right. \quad \frac{\frac{2}{P_1} + \frac{1}{R_1}}{2} = 1$$

Max f_1

Pr αF_1

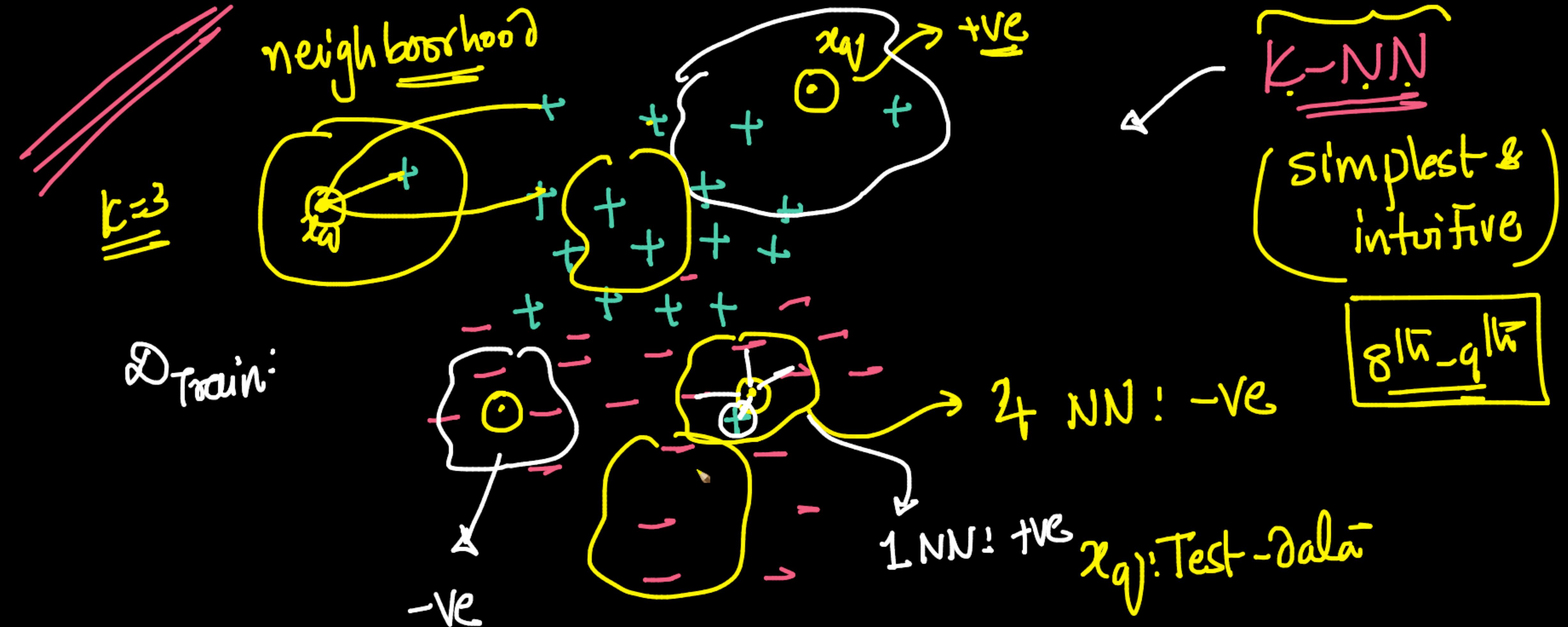
Re αF_1

Next-class: — ROC AUC
PR-Curves

→ Odds - ratio

→ imbalanced data

KNN



Inputs: $\mathcal{D}_{\text{Train}}$, x_q , $K = \boxed{\text{odd}} : 5$

$$\checkmark d = [$$

[for each x_i in $\mathcal{D}_{\text{Train}}$

- compute $\text{dist}(x_q, x_i)$

- append to d

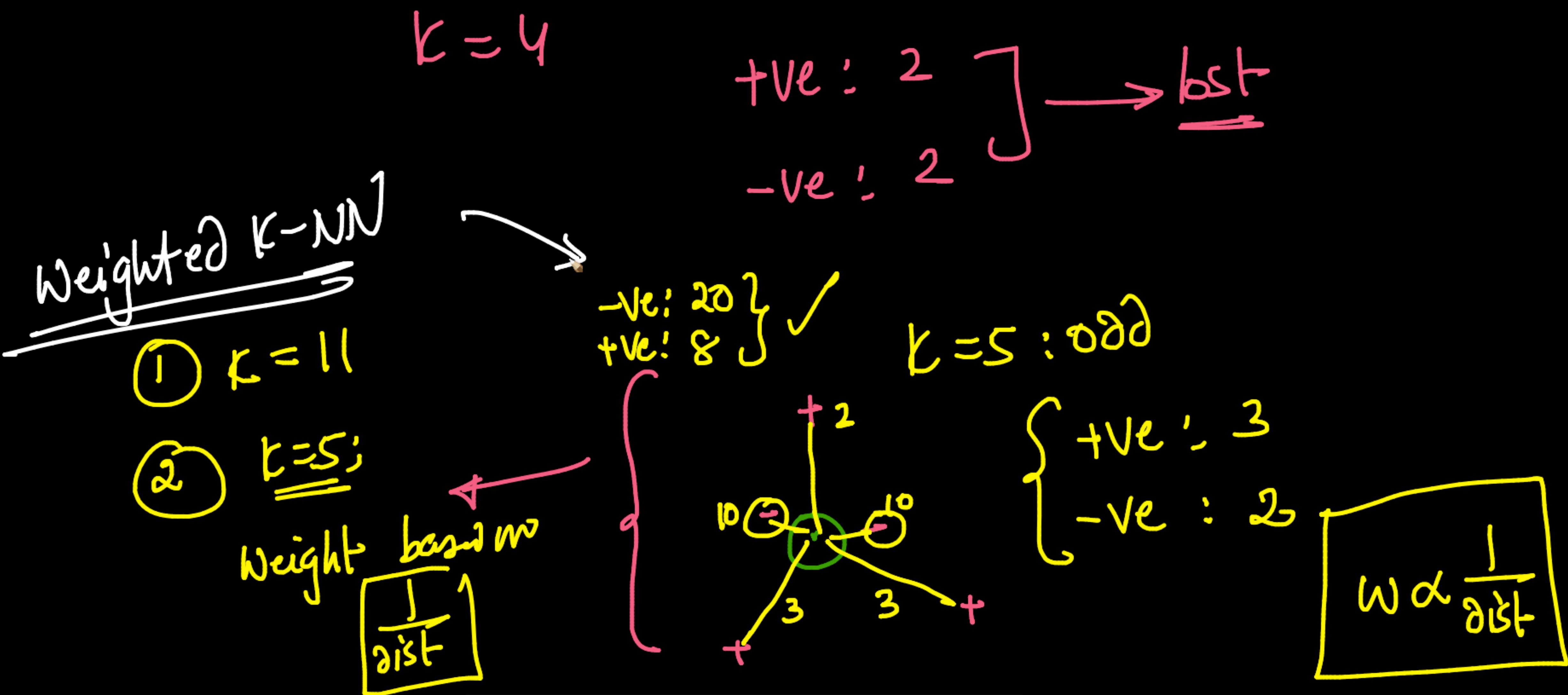
- **Sort** the $d[-]$

- pick \boxed{K} nearest neighbors

- Majority rule & return class label

$$\begin{aligned} & \text{dist}(x_q, x_i) \\ &= \sqrt{\sum_{j=1}^n (x_{qj} - x_{ij})^2} \\ &= \sqrt{\|x_q - x_i\|^2} \end{aligned}$$

+ve $\rightarrow \boxed{3}$
 -ve $\rightarrow \boxed{2}$



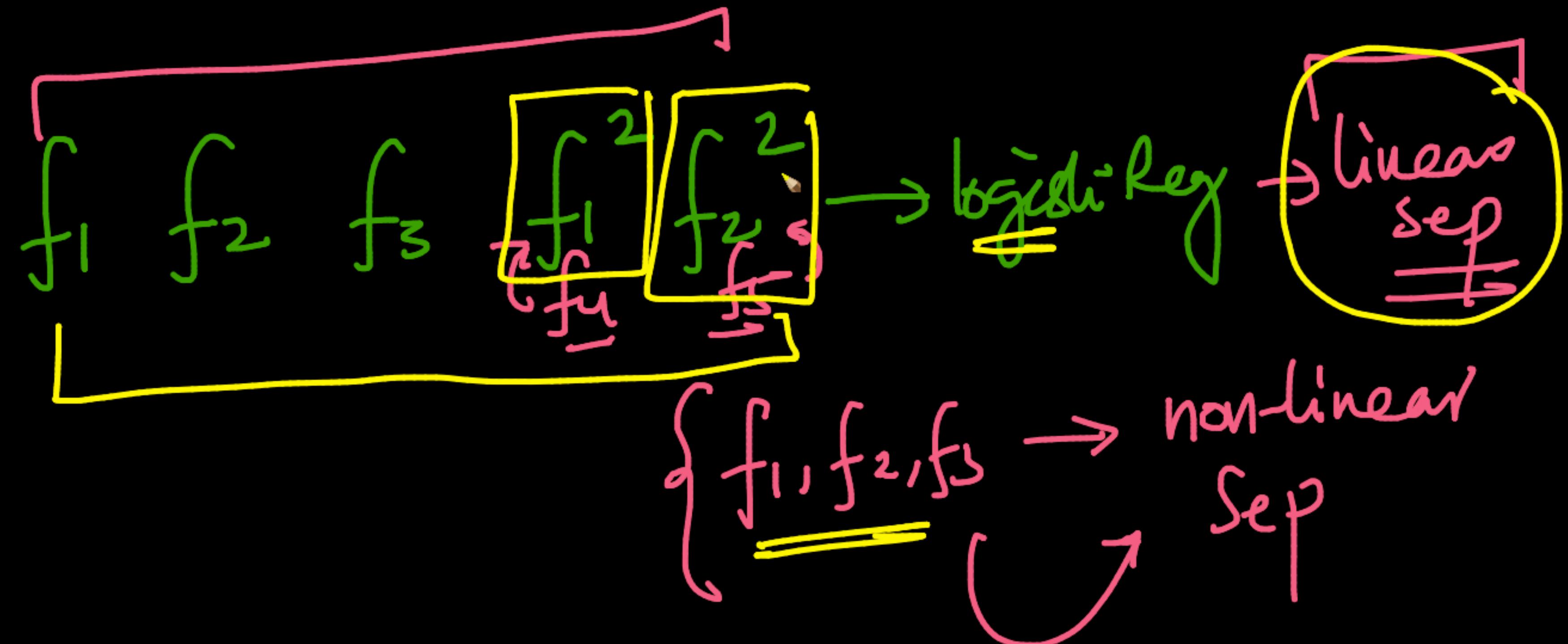
Assumption KNN:  Neighborhoods are homogenous most of the time

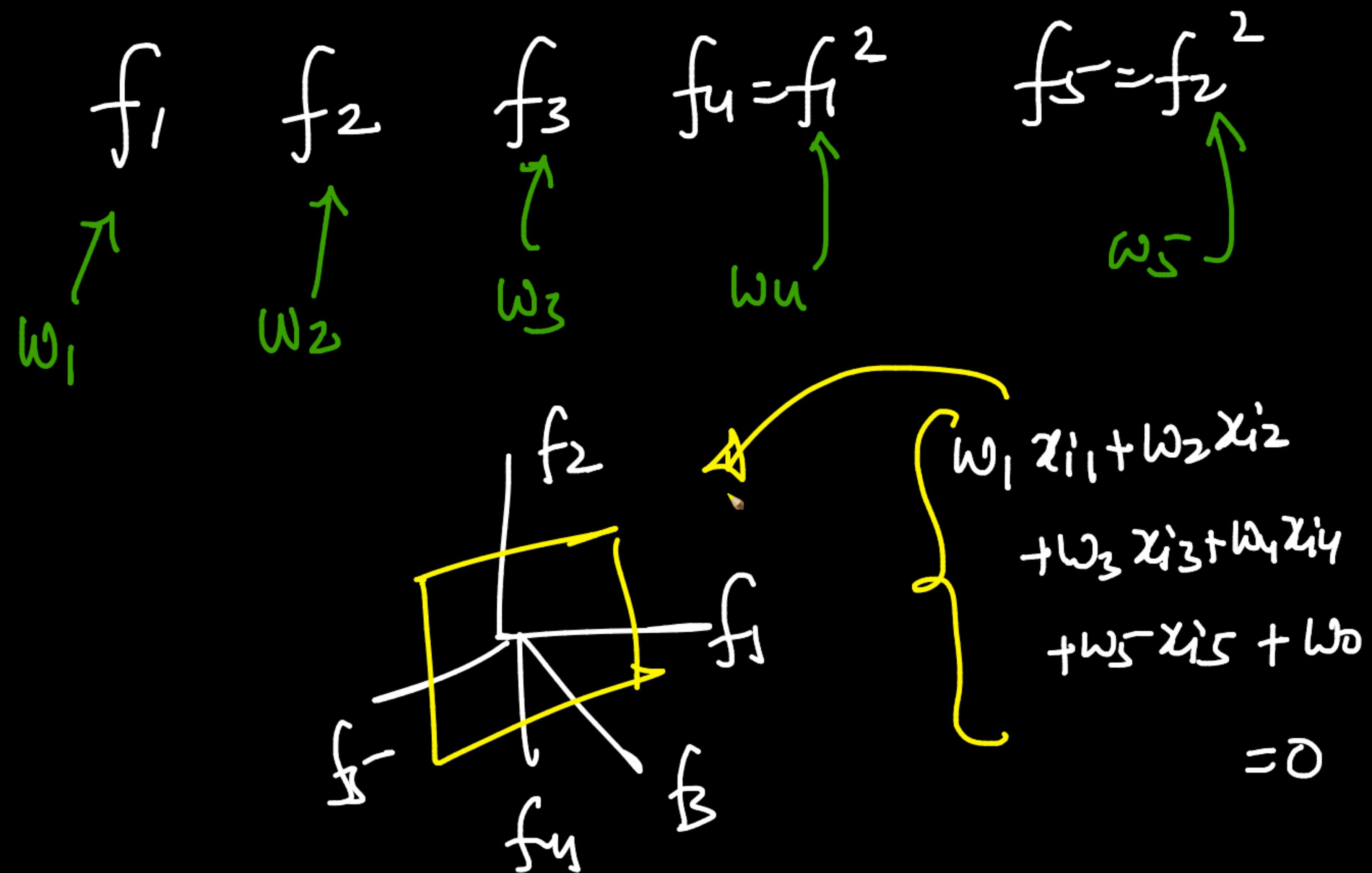
Assumption logistic Regr: → Linearity ~~Seperability~~

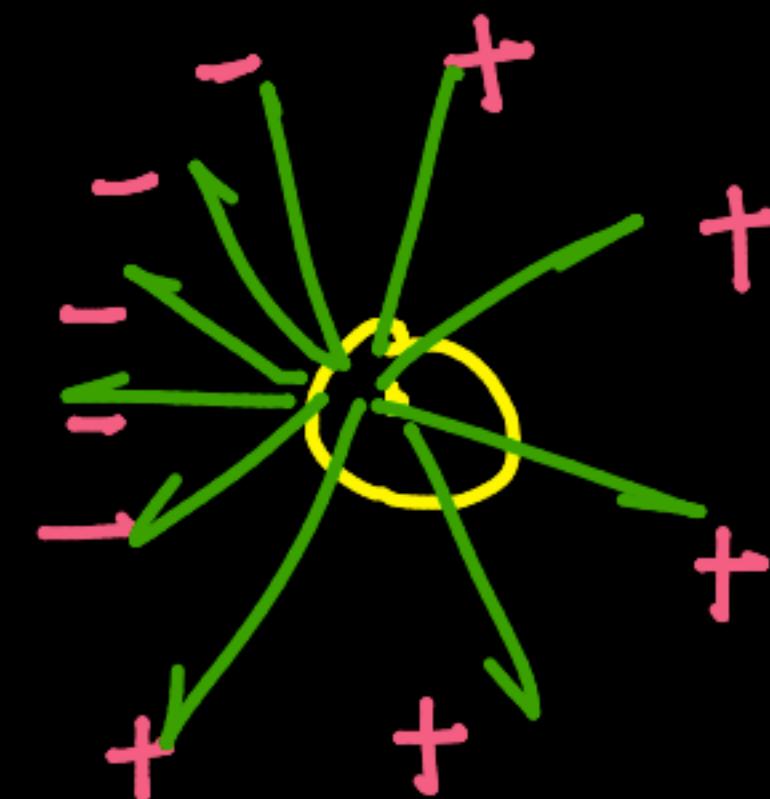
Good

$f_1 \ f_2 \ f_3$

\rightarrow logisticreg \rightarrow /



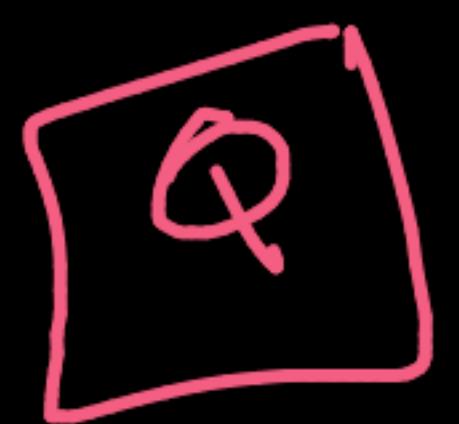




$k=5$

ties

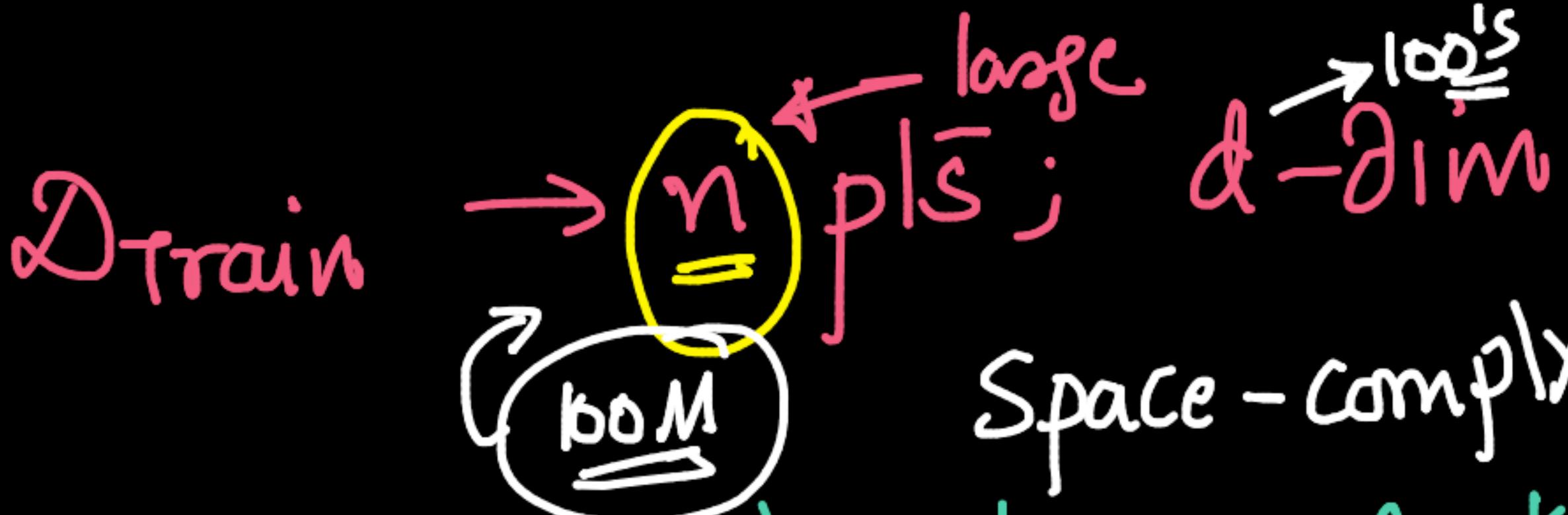
+ve : ↗
-ve : ↘



~~disadv~~

~~k-NN~~

(Test)



Space-complex

to run $k\text{-NN}$

$$\mathcal{O}(n \cdot d + \dots)$$

x_a

$x_a \rightarrow$

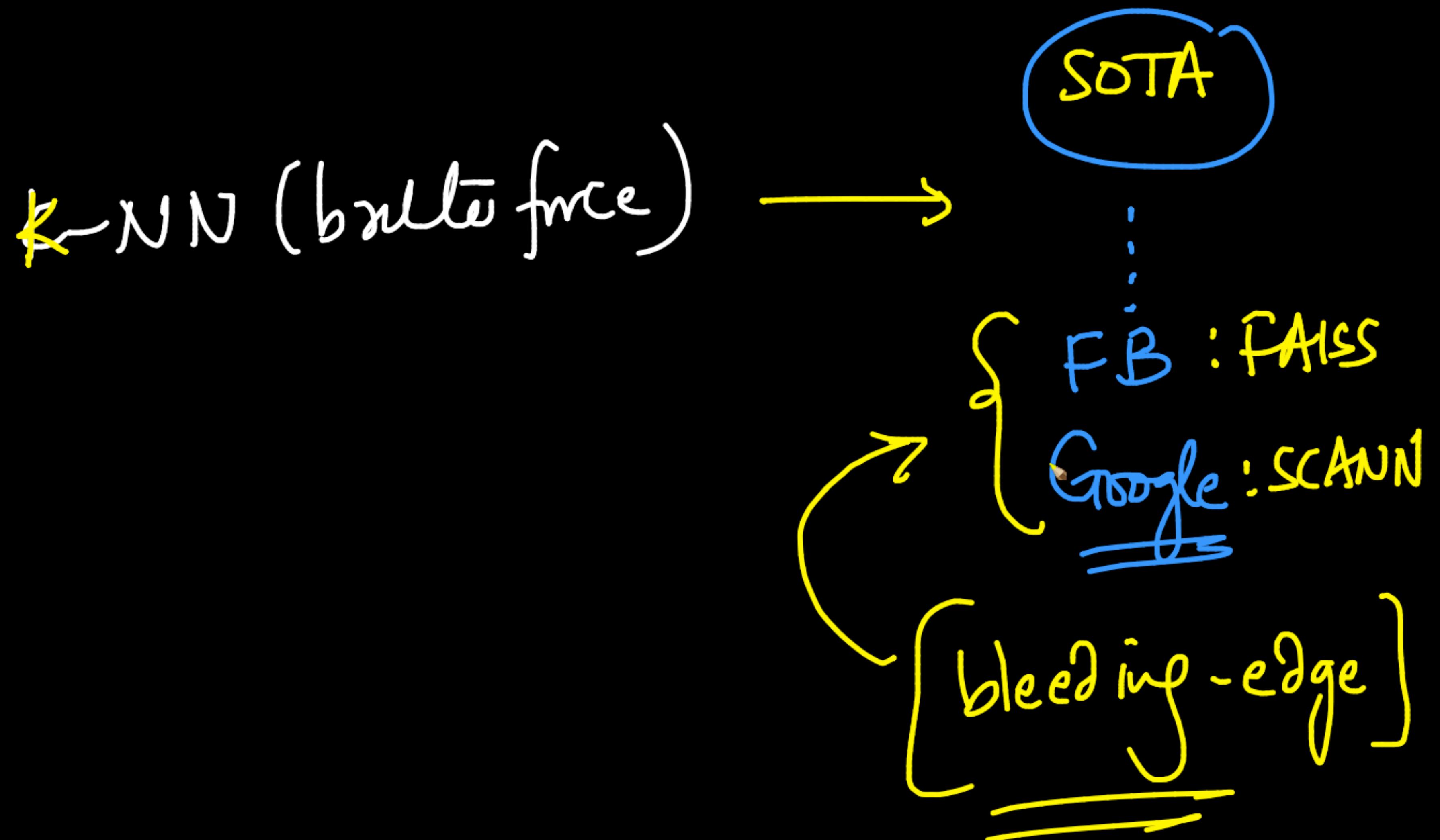
$\left\{ \begin{array}{l} - \text{dist}(x_i, x_a) \forall i \\ - n \text{ dist's} \cdot \text{Sort} \\ - k \cdot \text{top} \cdot \text{pls} \end{array} \right.$

BRUTE-FORCE

$$\mathcal{O}(n \cdot d + n \lg n + K)$$

$$\mathcal{O}(n \cdot (d + \lg n))$$

Small

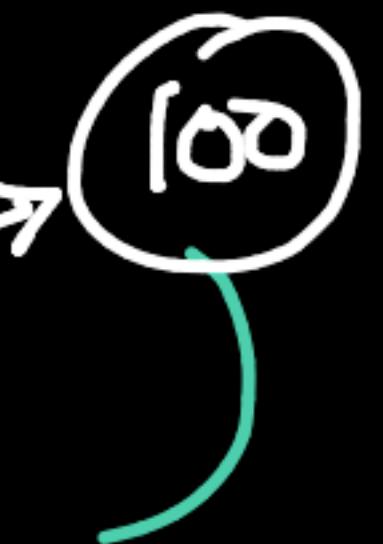


Q

✓ logistic reg: d-dim data n-train pls

$$\tau(w^T x_q + \underline{w_0})$$

d-mul

indep of n^2 

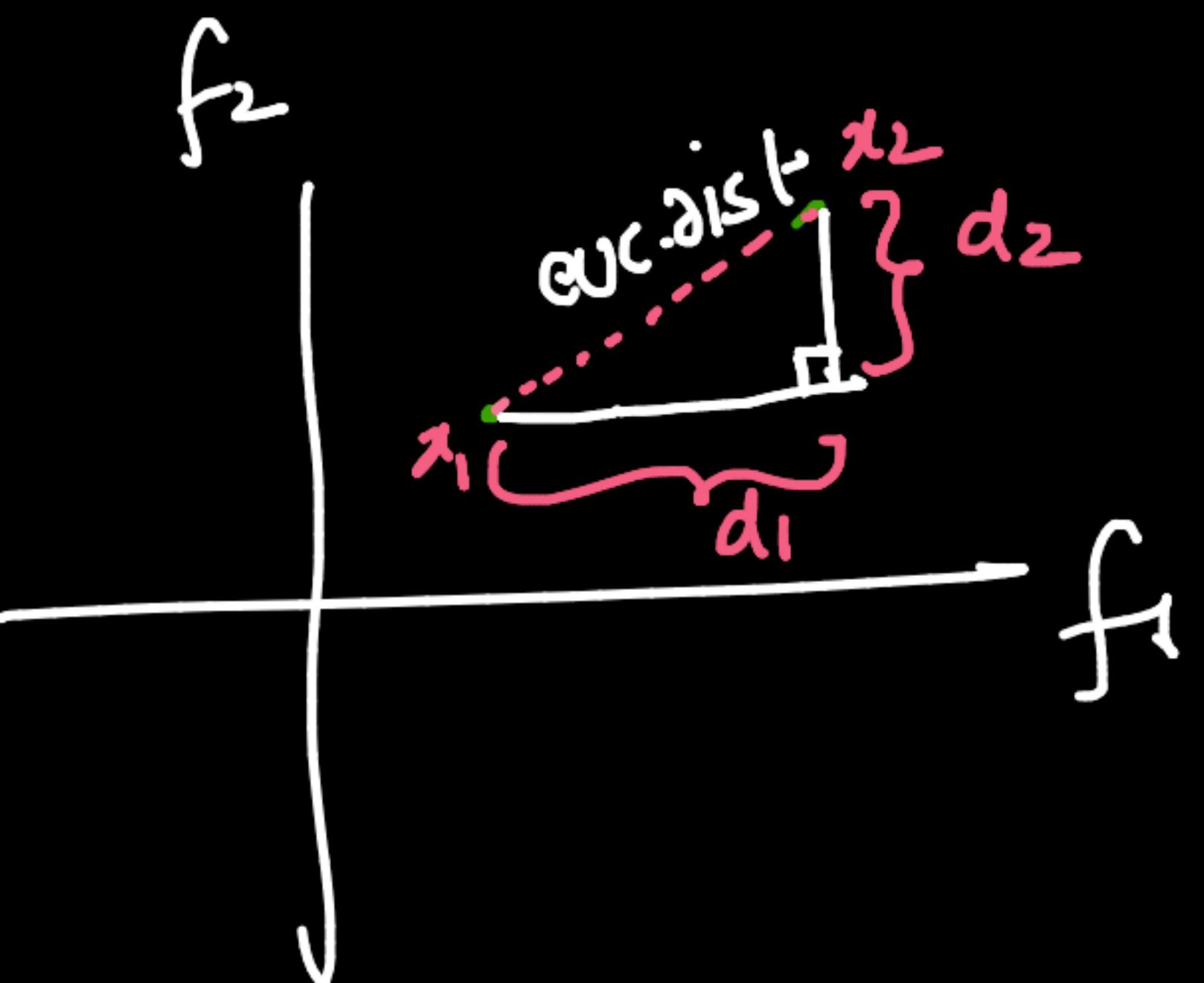
time-complex: $O(d)$

Space-Complex: $O(d)$

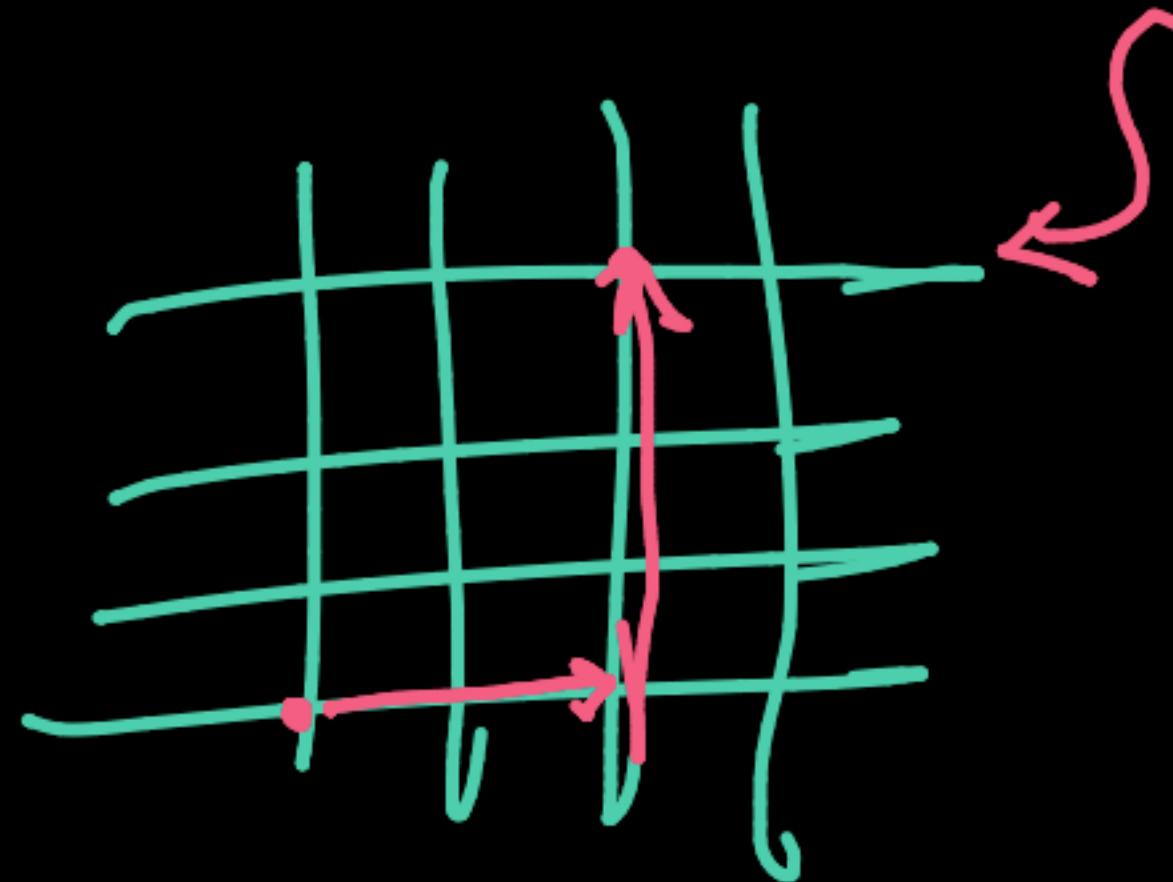
$$\text{euc. dist}(\mathbf{x}_1, \mathbf{x}_2) = \left[\sum_{j=1}^d (x_{1j} - x_{2j})^2 \right]^{1/2}$$

does not work
when d is high
(curse of dim) → Math (later) ...

$$\text{Manhattan dist}(x_1, x_2) = \left(\sum_{j=1}^d |x_{1j} - x_{2j}| \right)^{\frac{1}{2}}$$



$d_1 + d_2$: Manhattan dist



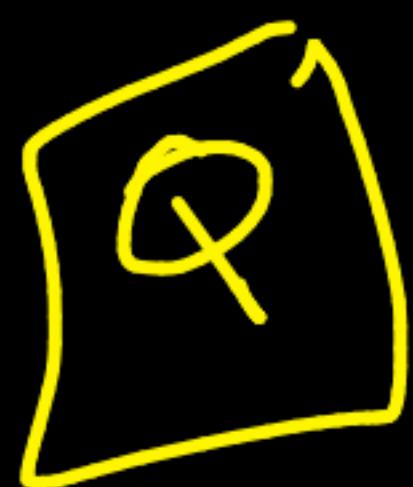
Minkowski dist (x_1, x_2, p)

generalization
of euc &
Manhattan
dist

$$\left(\sum_{j=1}^d |x_{1j} - x_{2j}|^p \right)^{1/p}$$

$p=2 \rightarrow$ euc.dist
 $p=1 \rightarrow$ Manhattan
dist

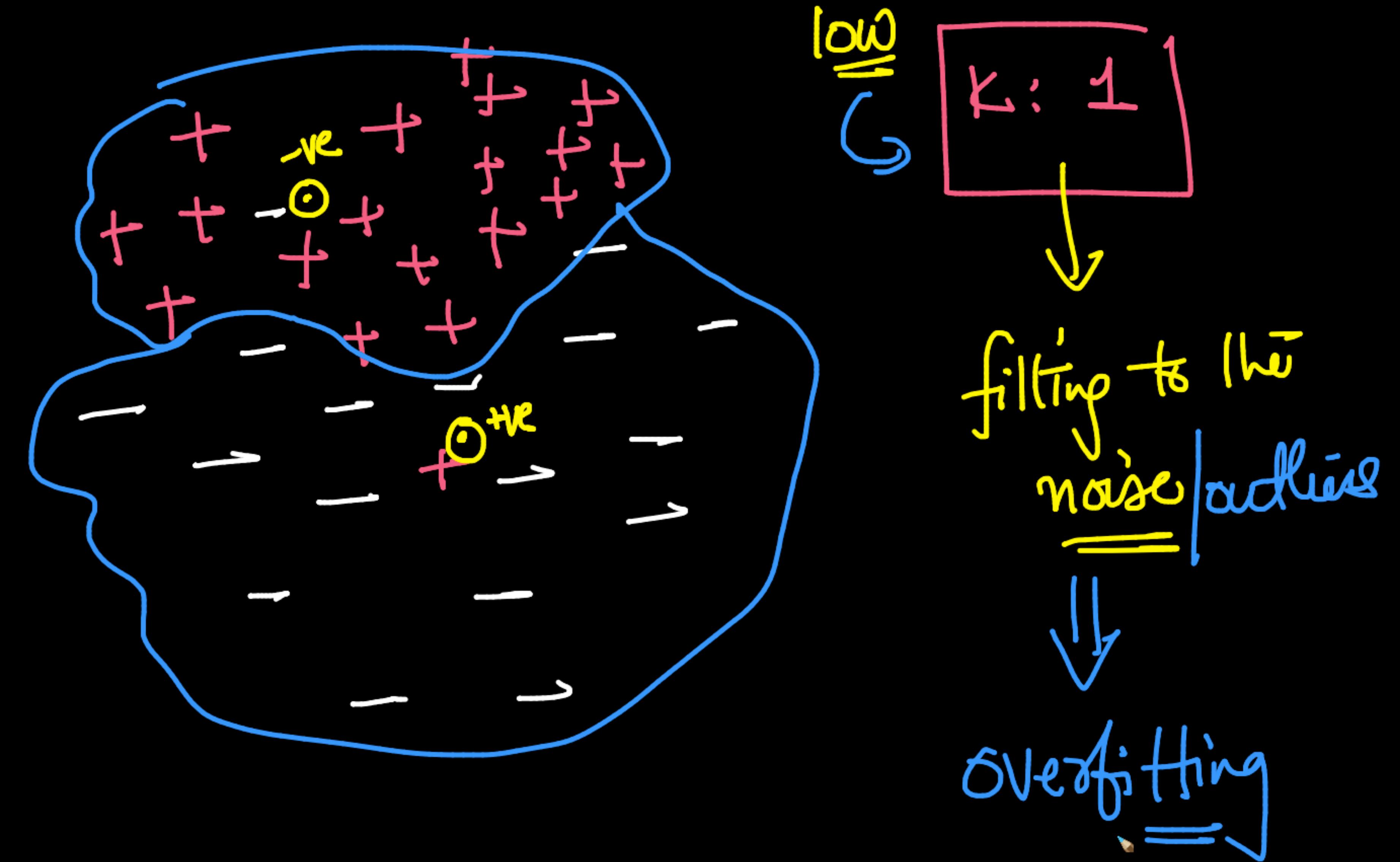
✓ { Cosine-Dist → Text Data

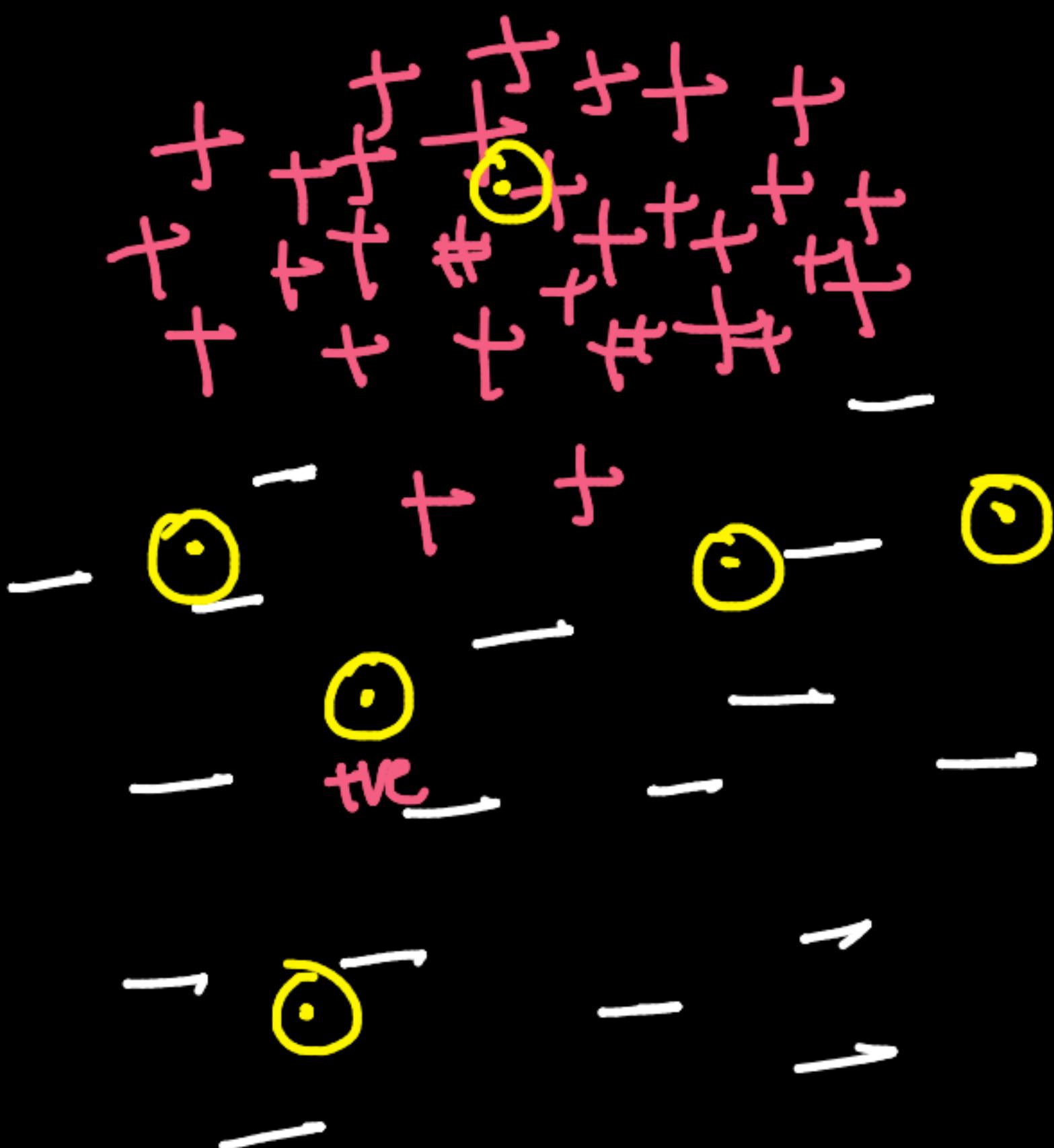


No training: $\mathcal{D}_{\text{train}}$

Test time ↗ \mathbf{x}_{qj}

bias - Variance
tradeoff 

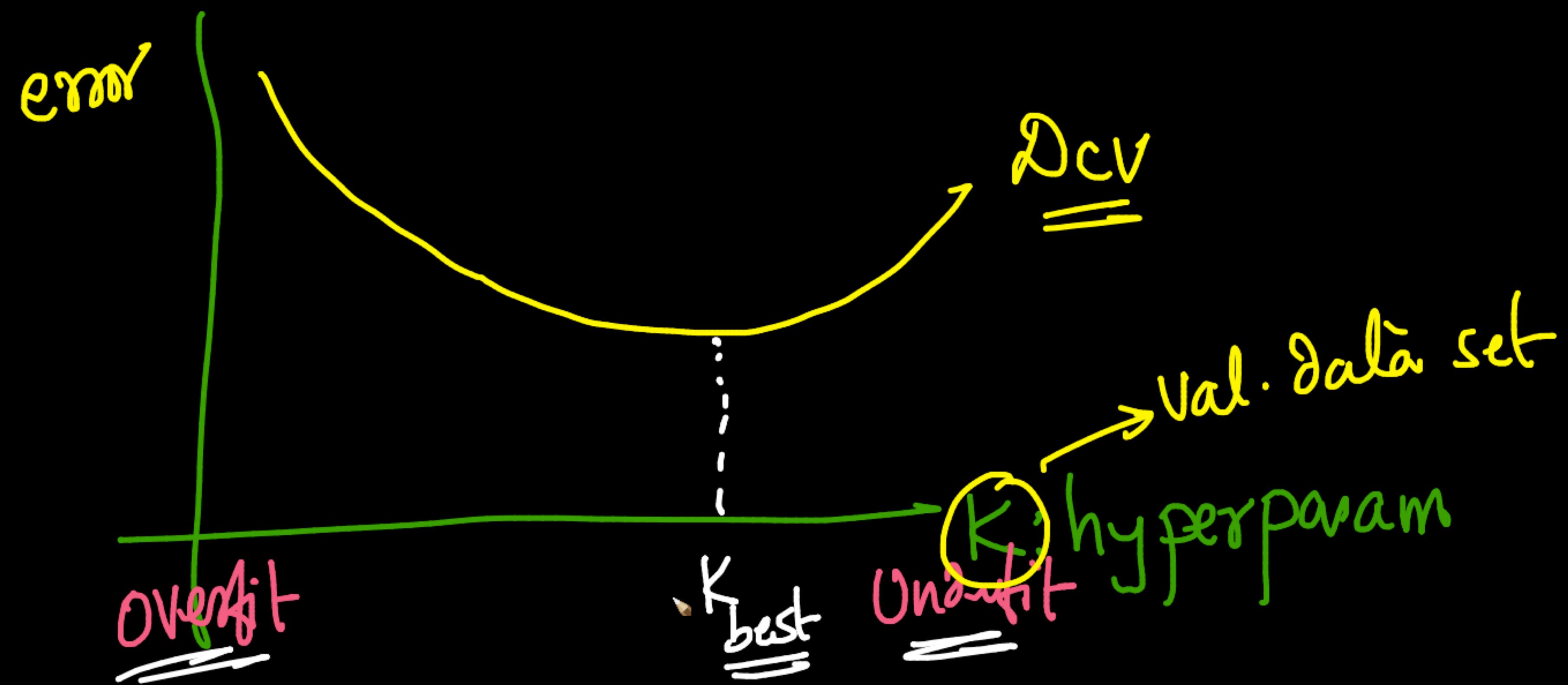


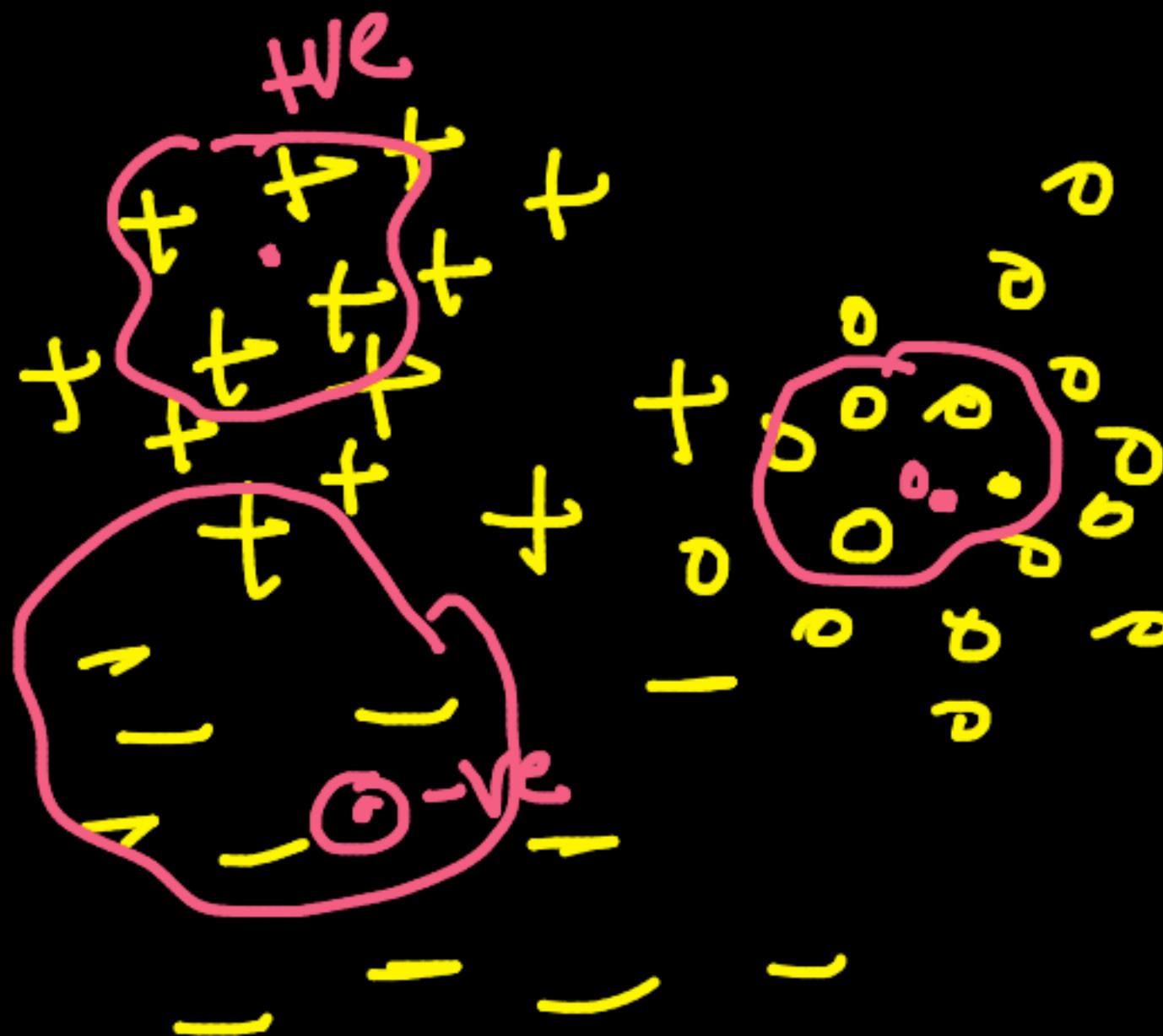
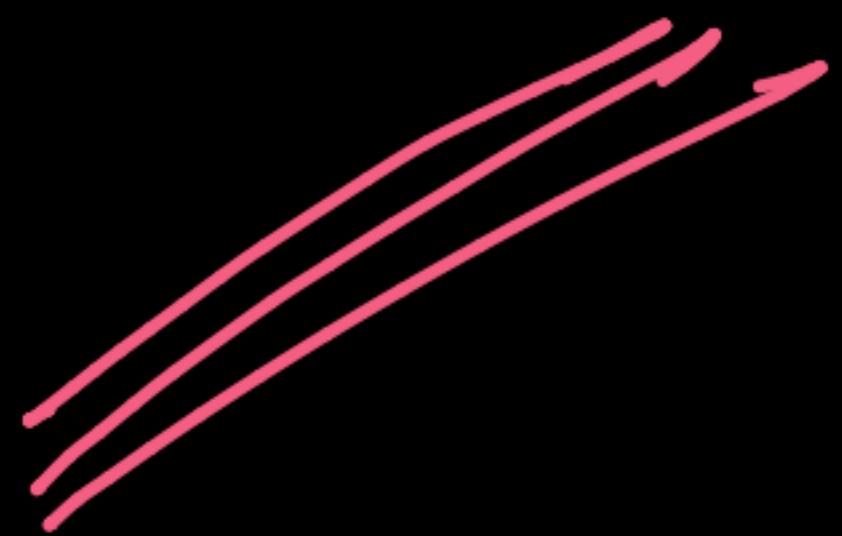
largest $K = n = \underline{10}$

$$\begin{cases} +ve: 61 \\ -ve: 40 \end{cases}$$

$y_{qj} = +ve$

Underfitting





multi-class
Classification



hyper-param

Majority rule
↳ class amongst K-NNs

K-NN

for other
algo

$$\mathcal{D}_{TR} = \mathcal{D}^0$$

$$\mathcal{D}^1$$

$$\mathcal{D}^2$$

:

day = 0

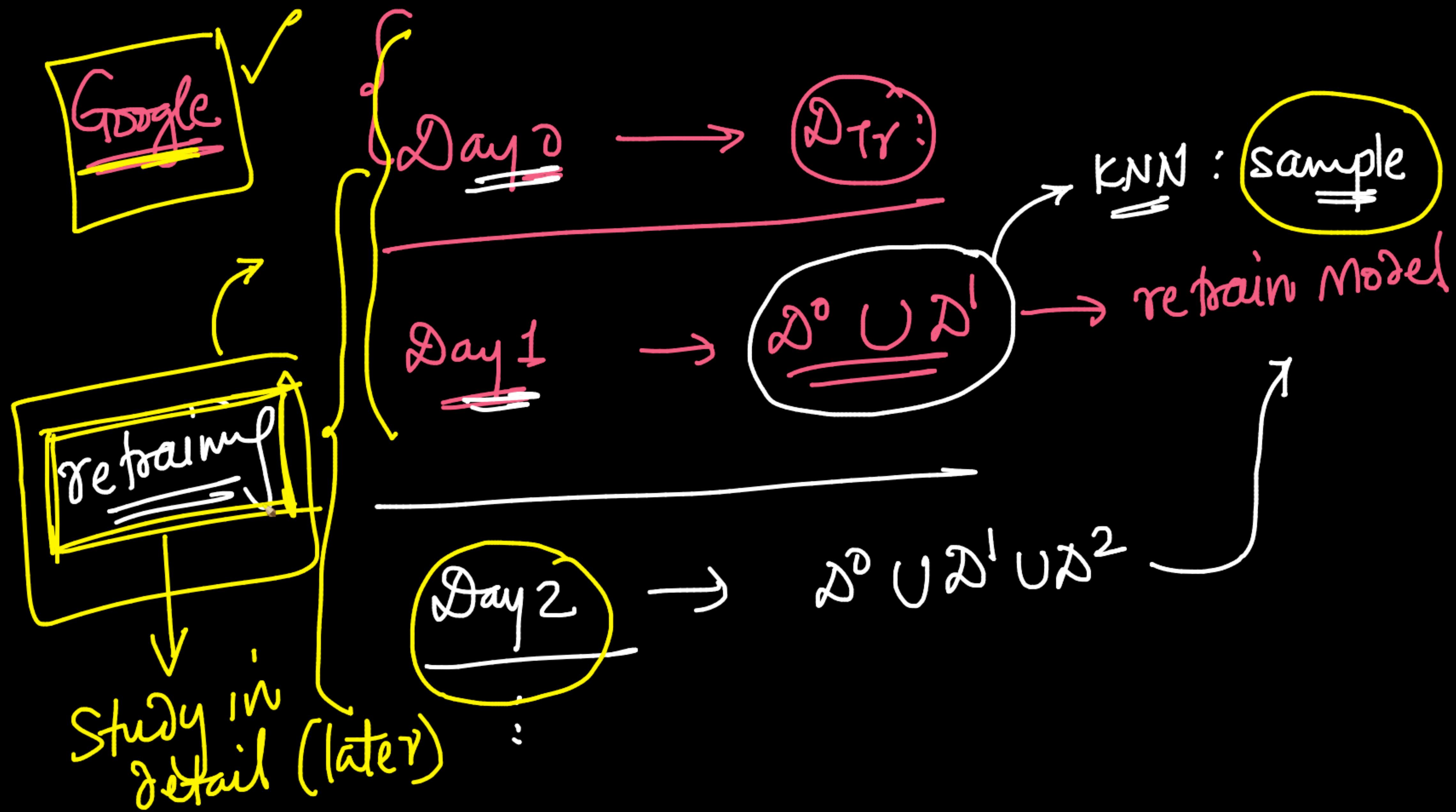
day = 1

day = 2

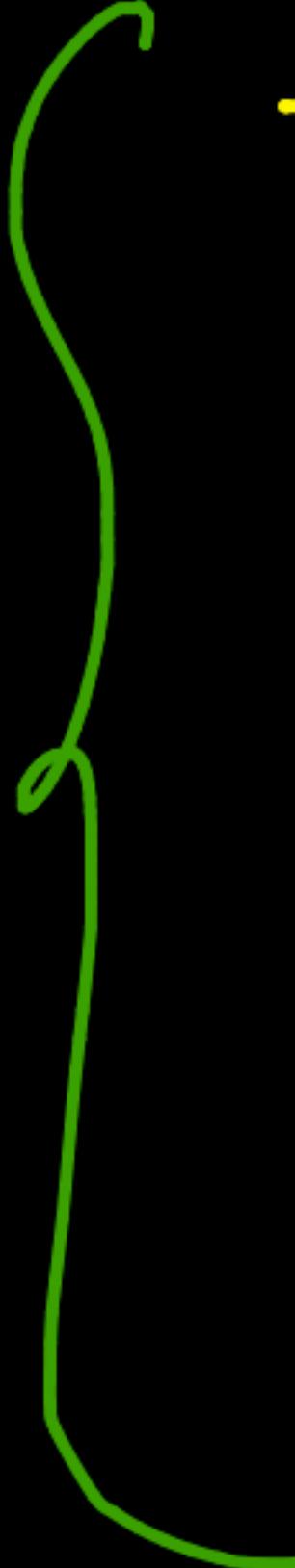
{ $\mathcal{D}^0 \cup \dots \cup \mathcal{D}^q$

$$\underline{\text{day}} = \underline{\mathcal{D}}$$

Model *



plan



-

k-NN (cont)

...

Misc

(next-class)

-

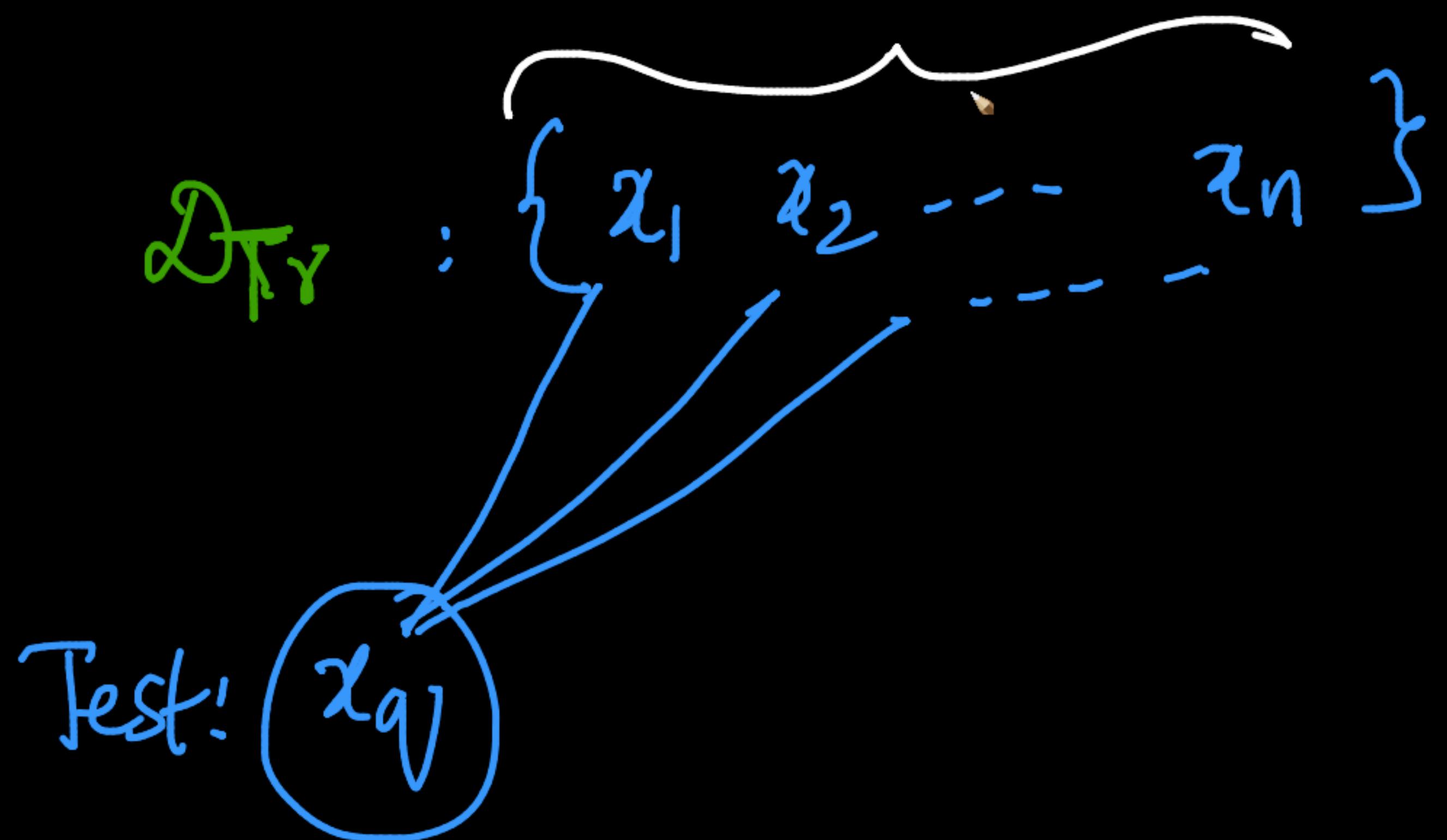
Odds ratio

-

AUC; ROC; PRC

:

!



logistic reg:

model: w_i^s

$$\frac{\partial L}{\partial w_j}$$

$$\min_{w_i} \sum_{i=1}^q \log \text{loss}_i$$

imbalanced data

$$L = \sum_{i=1}^q \log \text{loss}_i + \lambda \cdot L_2 \text{reg}$$

$$c_i = \begin{cases} 0.9 & \text{if } y_i = 1 \\ 0.1 & \text{if } y_i = 0 \end{cases}$$



The latest from Google Research

Announcing ScaNN: Efficient Vector Similarity Search

Tuesday, July 28, 2020

Posted by Philip Sun, Software Engineer, Google Research

Suppose one wants to search through a large dataset of literary works using queries that require an exact match of title, author, or other easily machine-indexable criteria. Such a task would be well suited for a relational database using a language such as SQL. However, if one wants to support more abstract queries, such as "Civil War poem," it is no longer possible to rely on naive similarity metrics such as the number of words in common between two phrases. For example, the query "science fiction" is more related to "future" than it is to "earth science" despite the former having zero, and the latter having one, word in common with the query.

Machine learning (ML) has greatly improved computers' abilities to understand language semantics and therefore answer these abstract queries. Modern ML models can transform inputs such as text and images into embeddings, high dimensional vectors trained such that more similar inputs cluster closer together. For a given query, we can therefore compute its embedding, and find the literary works whose embeddings are closest to the query's. In this manner, ML has transformed an abstract and previously difficult-to-specify task into a rigorous mathematical one.

However, a computational challenge remains: how do we efficiently search through a large dataset of embeddings to find the nearest dataset elements? This is known as vector similarity search and its high dimensionality makes pruning difficult.

In our [ICML 2020 paper](#), "[Accelerating Large-Scale Inference with Anisotropic Vector Quantization](#)," we address this problem by focusing on how to compress the dataset vectors to enable fast approximate distance computations, and propose a new compression technique that significantly boosts accuracy compared to prior works. This technique is utilized in our recently open-sourced [vector similarity search library](#) (ScaNN), and enables us to outperform other vector similarity search libraries by a factor of two, as measured on [ann-benchmarks.com](#).

The Importance of Vector Similarity Search

Embedding-based search is a technique that is effective at answering queries that rely on semantic understanding rather than simple indexable properties. In this technique, machine learning models

Search blog ...
 Labels
 Archive
 Feed
 Follow @googleai
Give us feedback in our [Product Forums](#).

