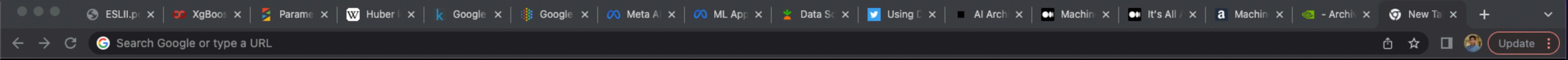


- light GBM - code → pseudo - residual
- Topics:
- GBDT : classification - loss
 - [- RF vs GBDT vs Cascading vs Stacking
 - Cascading & stacking → Questions
 - Kaggle vs real-world → ...
 - [- SVM
 - intution
 - Math (lots)



Google

Search Google or type a URL

Colaboratory My Drive Learning YouTube InterviewBit S...

GitHub Scaler Acad... Reduce the fil... InterviewBit Add shortcut

Keep your account safe and help prevent cyberattacks. Turn on 2-Step Verification

Photo by NASA Image Library

+ Code + Text

Reconnect ▾

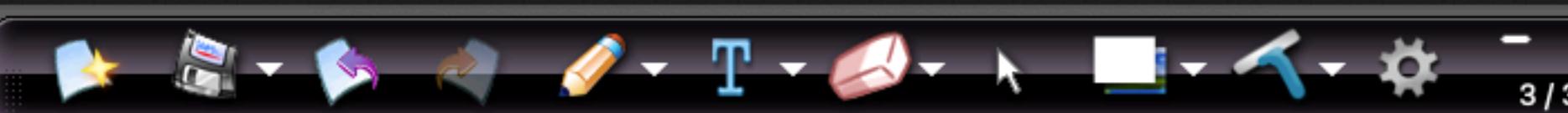
A small white icon representing user settings or account management.

1

```
# LightGBM: Sample code
import lightgbm as lgb
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
#Refer: https://lightgbm.readthedocs.io/en/latest/Parameters.html
```

```
gridParams = {  
    'learning_rate': [0.1, 0.5, 0.8],  
    'boosting_type' : ['gbdt'],  
    'objective' : ['multiclass'],  
    'max_depth' : [5,6,7,8],  
    'colsample_bytree' : [0.5,0.7],  
    'subsample' : [0.5,0.7],  
    'metric':['multi_error'],  
    'random_state' : [501]  
}
```

```
clf = lgb.LGBMRegressor(num_classes=20)
grid = RandomizedSearchCV(clf,gridParams,verbose=3,cv=3,n_jobs = -1,n_iter=10,
grid.fit(X_train,Y_train)
```



ESLII.pdf XgBoost.ipynb - Colaboratory Parameters – LightGBM 3.3.2 + colab.research.google.com/drive/1oOTcIDDRBX1mB3EZTsoYhv4S6Pk7TbjN#scrollTo=WAN5jxCLKCOh Reconnect + Code + Text

+ Code + Text

Reconnect

Code Editor

gridParams = {
 'learning_rate': [0.1, 0.5, 0.8],
 'boosting_type' : ['gbdt'],
 'objective' : ['multiclass'],
 'max_depth' : [5,6,7,8],
 'colsample_bytree' : [0.5,0.7],
 'subsample' : [0.5,0.7],
 'metric':['multi_error'],
 'random_state' : [501]
}

clf = lgb.LGBMRegressor(num_classes=20)
grid = RandomizedSearchCV(clf,gridParams,verbose=3,cv=3,n_jobs = -1,n_iter=10,)
grid.fit(X_train,Y_train)

Fitting 3 folds for each of 10 candidates, totalling 30 fits
/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_search.py:972: UserWarning: One or more of the test category=UserWarning,
RandomizedSearchCV(cv=3, estimator=LGBMRegressor(num_classes=20), n_jobs=-1,
param_distributions={'boosting_type': ['gbdt'],
 'colsample_bytree': [0.5, 0.7],
 'learning_rate': [0.1, 0.5, 0.8],
 'max_depth': [5, 6, 7, 8],
 'metric': ['multi_error']})

ESLII.pdf XgBoost.ipynb - Colaboratory Parameters – LightGBM 3.3.2.

lightgbm.readthedocs.io/en/latest/Parameters.html#metric-parameters

CONTENTS:

- Installation Guide
- Quick Start
- Python Quick Start
- Features
- Experiments

Parameters

- Parameters Format
- Core Parameters
- Learning Control Parameters
- IO Parameters
- Objective Parameters
- Metric Parameters
- Network Parameters
- GPU Parameters
- Others

Parameters Tuning

- CAPI
- Python API
- R API
- Distributed Learning Guide
- GPU Tutorial
- Advanced Topics
- FAQ
- Development Guide

v: latest

boosting_type 1/1

CONTENTS:

- requires labels in {0, 1}; see `cross-entropy` application for general probability labels in [0, 1]
- multi-class classification application
 - `multiclass`, `softmax` objective function, aliases: `softmax`
 - `multiclassova`, One-vs-All binary objective function, aliases: `multiclass_ova`, `ova`, `ovr`
 - `num_class` should be set as well
- cross-entropy application
 - `cross_entropy`, objective function for cross-entropy (with optional linear weights), aliases: `xentropy`
 - `cross_entropy_lambda`, alternative parameterization of cross-entropy, aliases: `xentlambda`
 - label is anything in interval [0, 1]
- ranking application
 - `lambdarank`, `lambdarank` objective. `label_gain` can be used to set the gain (weight) of `int` label and all values in `label` must be smaller than number of elements in `label_gain`
 - `rank_xendcg`, XE_NDCG_MART ranking objective function, aliases: `xendcg`, `xe_ndcg`, `xe_ndcg_mart`, `xendcg_mart`
 - `rank_xendcg` is faster than and achieves the similar performance as `lambdarank`
 - label should be `int` type, and larger number represents the higher relevance (e.g. 0:bad, 1:fair, 2:good, 3:perfect)
- boosting** `gbdt`, default = `gbdt`, type = enum, options `gbdt`, `rf`, `dart`, `goss`, aliases: `boosting_type`, `boost`
 - `gbdt`, traditional Gradient Boosting Decision Tree, aliases: `gbdt`
 - `rf`, Random Forest, aliases: `random_forest`
 - `dart`, Dropouts meet Multiple Additive Regression Trees
 - `goss`, Gradient-based One-Side Sampling
 - Note:** internally, LightGBM uses `gbdt` mode for the first `1 / learning_rate` iterations
- data** `train`, default = `""`, type = string, aliases: `train`, `train_data`, `train_data_file`, `data_filename`
 - path of training data, LightGBM will train from this data
 - Note:** can be used only in CLI version
- valid** `test`, default = `""`, type = string, aliases: `test`, `valid_data`, `valid_data_file`, `test_data`, `test_data_file`, `valid_filenames`
 - path(s) of validation/test data, LightGBM will output metrics for these data
 - support multiple validation data, separated by `,`
 - Note:** can be used only in CLI version
- num_iterations** `100`, default = `100`, type = int, aliases: `num_iteration`, `n_iter`, `num_tree`, `num_trees`, `num_round`, `num_rounds`, `nrounds`, `num_boost_round`, `n_estimators`
 - number of boosting iterations

ESLII.pdf XgBoost.ipynb - Colaboratory Parameters – LightGBM 3.3.2 + colab.research.google.com/drive/1oOTcIDDRBX1mB3EZTsoYhv4S6Pk7TbjN#scrollTo=WAN5jxCLKOOh Reconnect + Code + Text

gridParams = {
 'learning_rate': [0.1, 0.5, 0.8],
 'boosting_type' : ['gbdt'],
 'objective' : ['multiclass'],
 'max_depth' : [5,6,7,8]
 'colsample_bytree' : [0.5,0.7],
 'subsample' : [0.5,0.7],
 'metric':['multi_error'],
 'random_state' : [501]
}

clf = lgb.LGBMRegressor(num_classes=20)
grid = RandomizedSearchCV(clf,gridParams,verbose=3,cv=3,n_jobs = -1,n_iter=10,)
grid.fit(X_train,Y_train)

$P(y_i=1 | x_i)$
 $P(y_i=2 | x_i)$
:
 $P(y_i=20 | x_i)$

Fitting 3 folds for each of 10 candidates, totalling 30 fits
/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_search.py:972: UserWarning: One or more of the test category=UserWarning,
RandomizedSearchCV(cv=3, estimator=LGBMRegressor(num_classes=20), n_jobs=-1,
param_distributions={'boosting_type': ['gbdt'],
 'colsample_bytree': [0.5, 0.7],
 'learning_rate': [0.1, 0.5, 0.8],
 'max_depth': [5, 6, 7, 8],
 'metric': ['multi_error'],
 'objective': ['multiclass'],

ESLII.pdf XgBoost.ipynb - Colaboratory Parameters – LightGBM 3.3.2 + colab.research.google.com/drive/1oOTcIDDRBX1mB3EZTsoYhv4S6Pk7TbjN#scrollTo=WAN5jxCLKCOh Reconnect + Code + Text

gridParams = {
 'learning_rate': [0.1, 0.5, 0.8],
 'boosting_type' : ['gbdt'],
 'objective' : ['multiclass'],
 'max_depth' : [5,6,7,8],
 'colsample_bytree' : [0.5,0.7],
 'subsample' : [0.5,0.7],
 'metric':['multi_error'],
 'random_state' : [501]
}

clf = lgb.LGBMRegressor(num_classes=20)
grid = RandomizedSearchCV(clf,gridParams,verbose=3,cv=3,n_jobs = -1,n_iter=10,
grid.fit(X_train,Y_train)

3-fold CV
cores

Fitting 3 folds for each of 10 candidates, totalling 30 fits
/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_search.py:972: UserWarning: One or more of the test category=UserWarning,
RandomizedSearchCV(cv=3, estimator=LGBMRegressor(num_classes=20), n_jobs=-1,
param_distributions={'boosting_type': ['gbdt'],
 'colsample_bytree': [0.5, 0.7],
 'learning_rate': [0.1, 0.5, 0.8],
 'max_depth': [5, 6, 7, 8],
 'metric': ['multi_error'],
 'objective': ['multiclass'],



ESLII.pdf XgBoost.ipynb - Colaboratory Parameters – LightGBM 3.3.2 + colab.research.google.com/drive/1oOTcIDDRBX1mB3EZTsoYhv4S6Pk7TbjN#scrollTo=WAN5jxCLKCOh Reconnect + Code + Text

gridParams = {
 'learning_rate': [0.1, 0.5, 0.8],
 'boosting_type' : ['gbdt'],
 'objective' : ['multiclass'],
 'max_depth' : [5,6,7,8],
 'colsample_bytree' : [0.5,0.7],
 'subsample' : [0.5,0.7],
 'metric':['multi_error'],
 'random_state' : [501]
}

clf = lgb.LGBMRegressor(num_classes=20)
grid = RandomizedSearchCV(clf,gridParams,verbose=3,cv=3,n_jobs = -1,n_iter=10,)
grid.fit(X_train,Y_train)

Fitting 3 folds for each of 10 candidates, totalling 30 fits
/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_search.py:972: UserWarning: One or more of the test category=UserWarning,
RandomizedSearchCV(cv=3, estimator=LGBMRegressor(num_classes=20), n_jobs=-1,
param_distributions={'boosting_type': ['gbdt'],
 'colsample_bytree': [0.5, 0.7],
 'learning_rate': [0.1, 0.5, 0.8],
 'max_depth': [5, 6, 7, 8],
 'metric': ['multi_error'],
 'objective': ['multiclass'],

ESLII.pdf XgBoost.ipynb - Colaboratory Parameters – LightGBM 3.3.2 + colab.research.google.com/drive/1oOTcIDDRBX1mB3EZTsoYhv4S6Pk7TbjN#scrollTo=WAN5jxCLKCOh Reconnect + Code + Text

+ Code + Text

Reconnect

Code Editor

Run Cell

Search

Find

Replace

File

Cell

Kernel

Help

Play

Up

Down

Left

Right

Home

End

Copy

Cut

Paste

Clear

More

objective' : ['multiclass'],
'max_depth' : [5,6,7,8],
'colsample_bytree' : [0.5,0.7],
'subsample' : [0.5,0.7],
'metric':['multi_error'],
'random_state' : [501]
}

clf = lgb.LGBMRegressor(num_classes=20)
grid = RandomizedSearchCV(clf,gridParams,verbose=3,cv=3,n_jobs = -1,n_iter=10,)
grid.fit(X_train,Y_train)

Fitting 3 folds for each of 10 candidates, totalling 30 fits
/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_search.py:972: UserWarning: One or more of the test category=UserWarning,
RandomizedSearchCV(cv=3, estimator=LGBMRegressor(num_classes=20), n_jobs=-1,
param_distributions={'boosting_type': ['gbdt'],
'colsample_bytree': [0.5, 0.7],
'learning_rate': [0.1, 0.5, 0.8],
'max_depth': [5, 6, 7, 8],
'metric': ['multi_error'],
'objective': ['multiclass'],
'random_state': [501],
'subsample': [0.5, 0.7]},
verbose=1)

9/9

+ Code + Text

[]



LightGBM: Sample code
import lightgbm as lgb
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
[#Refer: https://lightgbm.readthedocs.io/en/latest/Parameters.html](https://lightgbm.readthedocs.io/en/latest/Parameters.html)

gridParams = {
 'learning_rate': [0.1, 0.5, 0.8], → 3
 'boosting_type' : ['gbdt'],
 'objective' : ['multiclass'],
 'max_depth' : [5,6,7,8], → 3
 'colsample_bytree' : [0.5,0.7], → 2
 'subsample' : [0.5,0.7], → 2
 'metric':['multi_error'],
 'random_state' : [501]
}

clf = lgb.LGBMRegressor(num_classes=20)
grid = RandomizedSearchCV(clf,gridParams,verbose=3,cv=3,n_jobs = -1,n_iter=10,)
grid.fit(X_train,Y_train)

$$9 \times 4 = \underline{\underline{36}}$$

10

randomly

Fitting 3 folds for each of 10 candidates, totalling 30 fits

/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_search.py:972: UserWarning: One or more of the test category=UserWarning,

ESLII.pdf XgBoost.ipynb - Colaboratory Parameters – LightGBM 3.3.2 + colab.research.google.com/drive/1oOTcIDDRBX1mB3EZTsoYhv4S6Pk7TbjN#scrollTo=WAN5jxCLKCOh Reconnect + : + Code + Text

[] value value value value = value = value = [11, 44, 23, 25, 41, 73, 21, 227, 31, 39, 0, 0
188, 0, 0, 0, 0, 14, 0, 30]

{x}

[] # Xgboost

```
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
from sklearn.model_selection import StratifiedKFold

import datetime as dt

params = {
    'learning_rate': [0.1, 0.5, 0.8], — 3
    'subsample': [0.6, 0.8, 1.0], — 3
    'colsample_bytree': [0.6, 0.8, 1.0], — 3
    'max_depth': [3, 4, 5] — 3
}
xgb = XGBClassifier(n_estimators=100, objective='multi:softmax', num_class=20, silent=True)
```

3⁴ = 81

[] folds = 3

11 / 11

Light GBM: (2017 research paper)

G :

ESLII.pdf XgBoost.ipynb - Colaboratory Parameters – LightGBM 3.3.2 LightGBM: A Highly Efficient Gradient Boosting Decision Tree

proceedings.neurips.cc/paper/2017/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf

3 / 9 | - 175% + | ☰ ⌂

LightGBM: A Highly Efficient Gradient Boosting Decision Tree

```

Input:  $m$ : feature dimension
 $\text{nodeSet} \leftarrow \{0\}$  ▷ tree nodes in current level
 $\text{rowSet} \leftarrow \{\{0, 1, 2, \dots\}\}$  ▷ data indices in tree nodes
for  $i = 1$  to  $d$  do
    for  $\text{node}$  in  $\text{nodeSet}$  do
         $\text{usedRows} \leftarrow \text{rowSet}[\text{node}]$ 
        for  $k = 1$  to  $m$  do
             $H \leftarrow \text{new Histogram()}$ 
            ▷ Build histogram
            for  $j$  in  $\text{usedRows}$  do
                 $\text{bin} \leftarrow I.f[k][j].\text{bin}$ 
                 $H[\text{bin}].y \leftarrow H[\text{bin}].y + I.y[j]$ 
                 $H[\text{bin}].n \leftarrow H[\text{bin}].n + 1$ 
            Find the best split on histogram  $H$ .
            ...
        Update  $\text{rowSet}$  and  $\text{nodeSet}$  according to the best split points.
        ...
    
```

```

Input:  $a$ : sampling ratio of large gradient data
Input:  $b$ : sampling ratio of small gradient data
Input:  $loss$ : loss function,  $L$ : weak learner
 $\text{models} \leftarrow \{\}$ ,  $\text{fact} \leftarrow \frac{1-a}{b}$ 
 $\text{topN} \leftarrow a \times \text{len}(I)$ ,  $\text{randN} \leftarrow b \times \text{len}(I)$ 
for  $i = 1$  to  $d$  do
     $\text{preds} \leftarrow \text{models.predict}(I)$ 
     $g \leftarrow loss(I, \text{preds})$ ,  $w \leftarrow \{1, 1, \dots\}$ 
     $\text{sorted} \leftarrow \text{GetSortedIndices}(\text{abs}(g))$ 
     $\text{topSet} \leftarrow \text{sorted}[1:\text{topN}]$ 
     $\text{randSet} \leftarrow \text{RandomPick}(\text{sorted}[\text{topN}:\text{len}(I)], \text{randN})$ 
     $\text{usedSet} \leftarrow \text{topSet} + \text{randSet}$ 
     $w[\text{randSet}] \times = \text{fact}$  ▷ Assign weight  $fact$  to the small gradient data.
     $\text{newModel} \leftarrow L(I[\text{usedSet}], -g[\text{usedSet}], w[\text{usedSet}])$ 
     $\text{models.append(newModel)}$ 

```

- $\frac{\partial L}{\partial F_M(x)}$

Tim

more pls
large residual

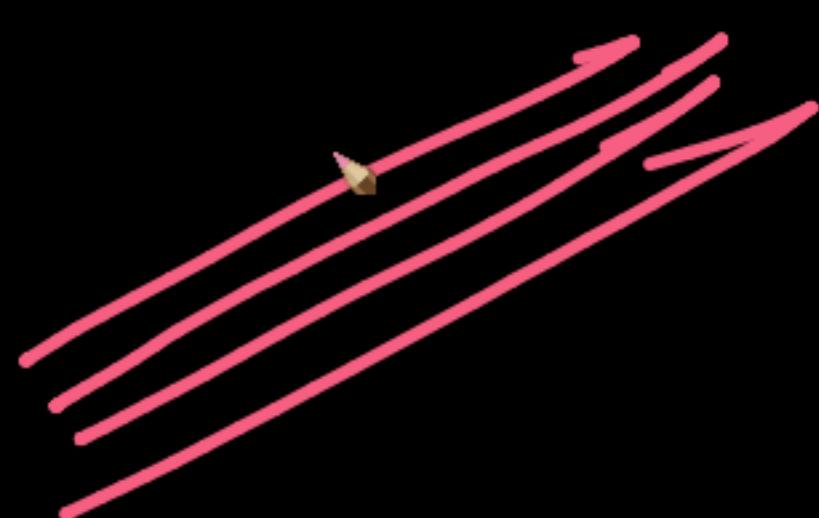
3 Gradient-based One-Side Sampling

In this section, we propose a novel sampling method for GBDT that can achieve a good balance between reducing the number of data instances and keeping the accuracy for learned decision trees.

3.1 Algorithm Description

In AdaBoost, the sample weight serves as a good indicator for the importance of data instances. However, in GBDT, there are no native sample weights, and thus the sampling methods proposed for AdaBoost cannot be directly applied. Fortunately, we notice that the gradient for each data instance in GBDT provides us with useful information for data sampling. That is, if an instance is associated with a small gradient, the training error for this instance is small and it is already well-trained. A straightforward idea is to discard those data instances with small gradients. However, the data distribution will be changed by doing so, which will hurt the accuracy of the learned model. To avoid this problem, we propose a new method called Gradient-based One-Side Sampling (GOSS).

3



ESLII.pdf XgBoost.ipynb - Colaboratory | Parameters – LightGBM 3.3.2 | Distributed Learning Guide – L x +

hastie.su.domains/Papers/ESLII.pdf

379 / 764 | - 200% + | ☰ 🔍

ESLII.pdf

360 10. Boosting and Additive Trees

TABLE 10.2. Gradients for commonly used loss functions.

Setting	Loss Function	$-\partial L(y_i, f(x_i))/\partial f(x_i)$
Regression	$\frac{1}{2}[y_i - f(x_i)]^2$	$y_i - f(x_i)$
Regression	$ y_i - f(x_i) $	$\text{sign}[y_i - f(x_i)]$
Regression	Huber	$y_i - f(x_i)$ for $ y_i - f(x_i) \leq \delta_m$ $\delta_m \text{sign}[y_i - f(x_i)]$ for $ y_i - f(x_i) > \delta_m$ where $\delta_m = \alpha\text{th-quantile}\{ y_i - f(x_i) \}$
Classification	Deviance	$k\text{th component: } I(y_i = g_k) - p_k(x_i)$

$-\frac{\partial L}{\partial F(x)}$

$\log \text{loss}$

residual

boosting procedure, and top-down decision tree induction, are themselves approximation procedures. After constructing the tree (10.37), the corresponding constants in each region are given by (10.30).

Table 10.2 summarizes the gradients for commonly used loss functions. For squared error loss, the negative gradient is just the ordinary residual $-g_{im} = y_i - f_{m-1}(x_i)$, so that (10.37) on its own is equivalent standard least squares boosting. With absolute error loss, the negative gradient is

tree to the progression the

Loss: log-loss

binary classfn: -

To show:

$$\frac{-\partial L}{\partial f_k(x_i)}$$

= pseudo residual \approx residual
 $\text{diff}(y_i, p_i)$

$$f_k(x_i) = \boxed{p_i} = P(y_i=1 | x_i)$$

y_i = actual class label or 1

✓
$$\frac{\partial L(y_i, p_i)}{\partial p_i} = -\frac{y_i}{p_i} + (1-y_i) \left(+1 \right) \frac{1}{1-p_i}$$

$$L = \text{Log-loss}(y_i, p_i) = y_i \log p_i + (1-y_i) \log(1-p_i)$$

$$-\frac{\partial L}{\partial p_i} = \frac{1-y_i}{1-p_i} - \frac{y_i}{p_i}$$

Deviance \rightarrow
 $p_i - y_i$

$$= \frac{p_i(1-y_i) - y_i(1-p_i)}{p_i(1-p_i)}$$

$$-\frac{\partial L}{\partial p_i} = \frac{p_i - y_i}{p_i(1-p_i)} - \frac{y_i}{p_i} + \frac{y_i}{1-p_i}$$

\checkmark
 $\frac{p_i - y_i}{p_i(1-p_i)}$

y_i : actual label 0 or 1

$$\hat{p}_i = p(y_i=1 | x_i)$$

pseudo residual = $\frac{\hat{p}_i - y_i}{(\hat{p}_i)(1-\hat{p}_i)}$

$$\text{residual} = \hat{p}_i - y_i$$
$$\text{diff}(\hat{p}_i, y_i)$$

SQ-loss:

Pseudo residual = $\sqrt{2(\hat{y}_i - \bar{y}_i)}$

Ques

pseudo residual \times

good estimate of
diff loss fn

✓ pseudo-residual
of any diff-
loss fn

\approx sensible
residual
 $\underline{\text{diff}}(y_i, \hat{y}_i)$

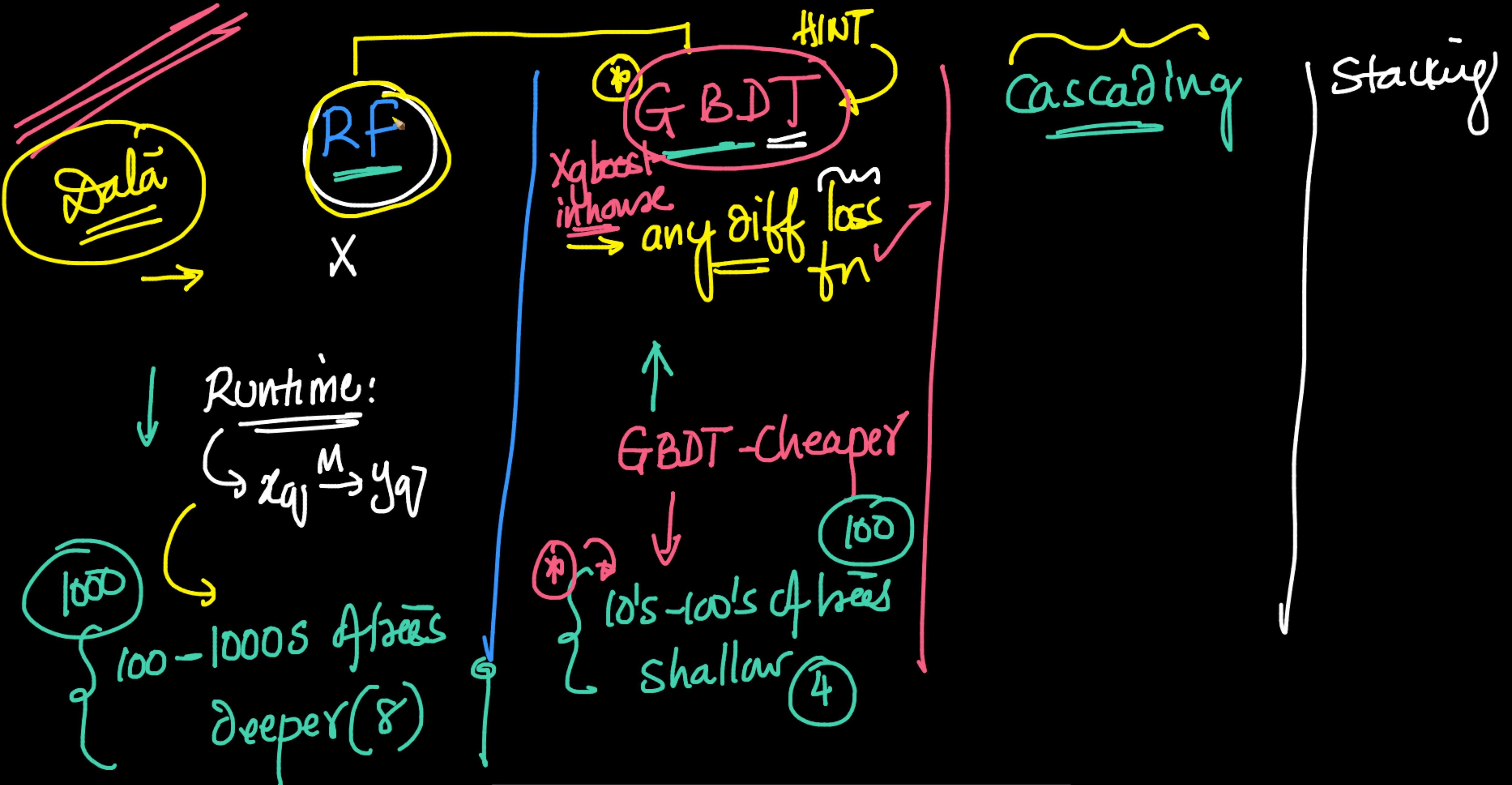




Derivative
~~log-odds~~ → logit formulation

$y_i - \hat{p}_i$

logistic reg: → linear reg on ~~log-odds~~



Deployment of GBDT:

$$F_M(x) = \underbrace{1 \cdot f_0(x)}_{\text{regy}} + \alpha_1 \underbrace{f_1(x)}_{\gamma^{1.6} \rightarrow 1.2} + \alpha_2 \underbrace{f_2(x)}_{\gamma^{0.2} \rightarrow 1.8} + \dots + \alpha_M \underbrace{f_M(x)}_{\gamma}$$

~~regy~~

x_i → each $\underline{\underline{f_i(x)}}$

$f_i(x)$: shallow DT

~~reg:~~

$$x_i : \begin{array}{|c|c|c|c|} \hline & 1 & 2 & 1 & 0 \\ \hline \text{dim} & \swarrow & \searrow & \swarrow & \searrow \\ \end{array}$$
$$y_i = \boxed{36}$$

errors

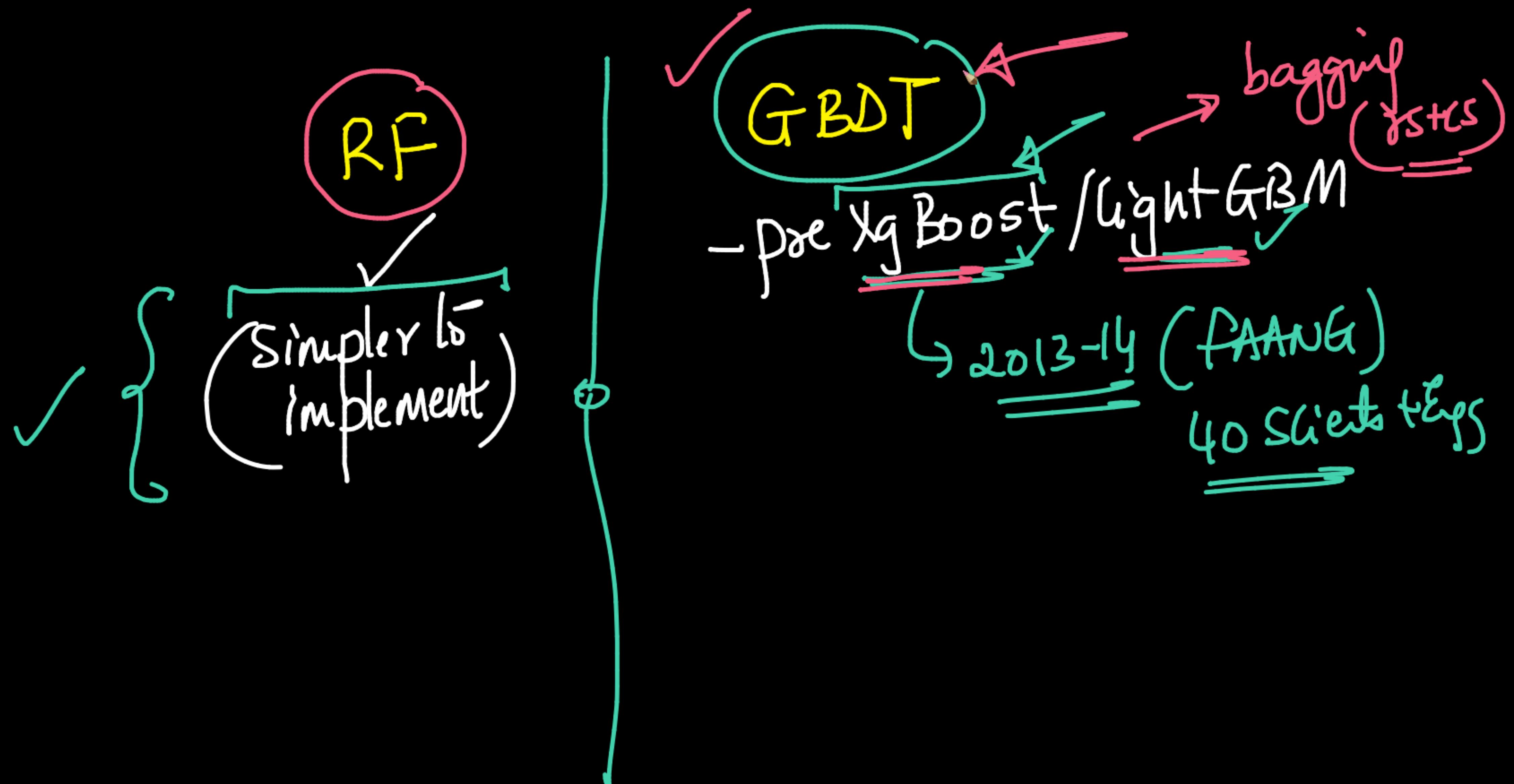
$$f_0(x_i) = \circled{12} \rightarrow 24$$

$$f_1(x_i) = \circled{18} \rightarrow 6$$

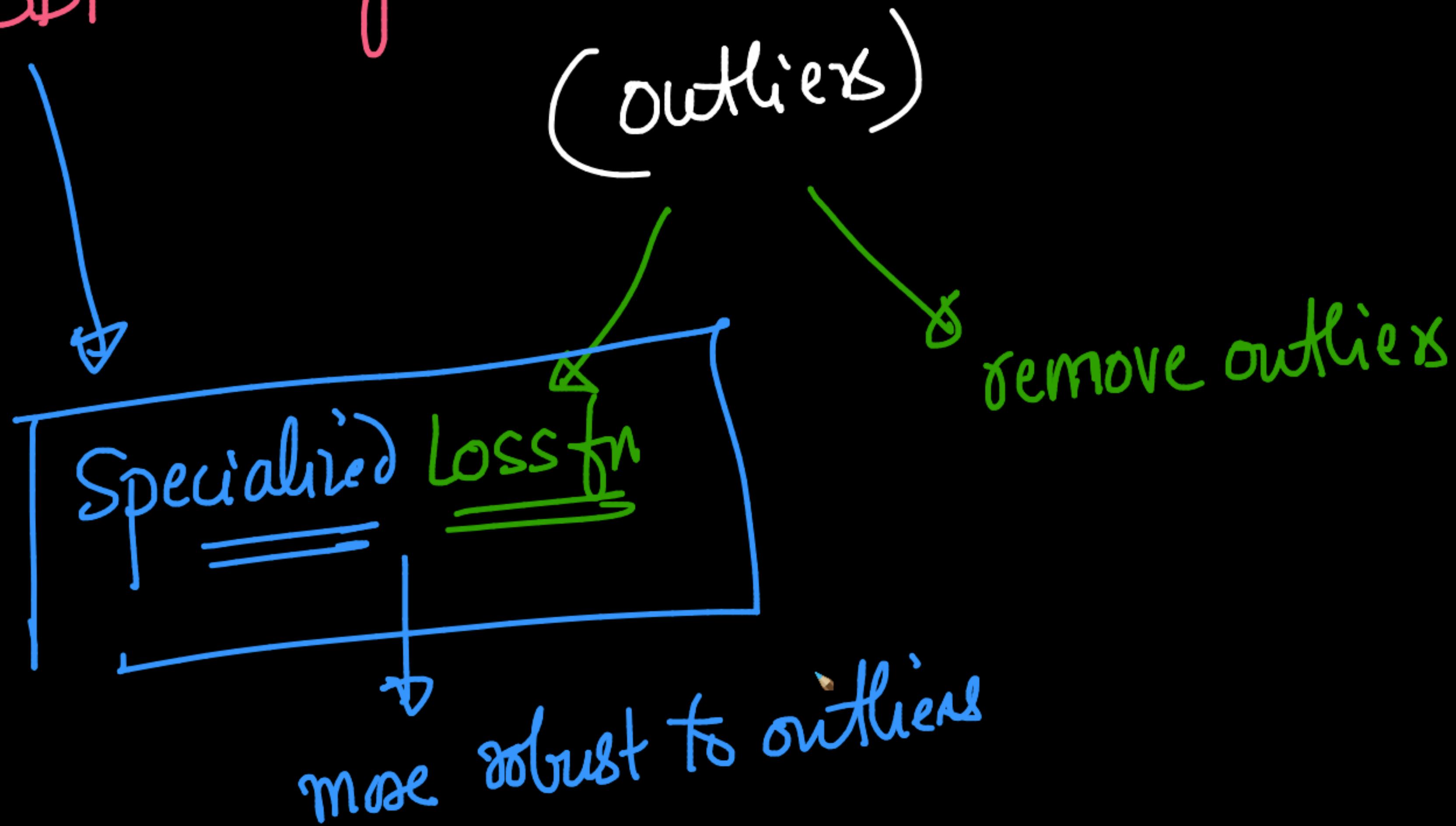
$$f_2(x_i) = \circled{4} \rightarrow 2$$

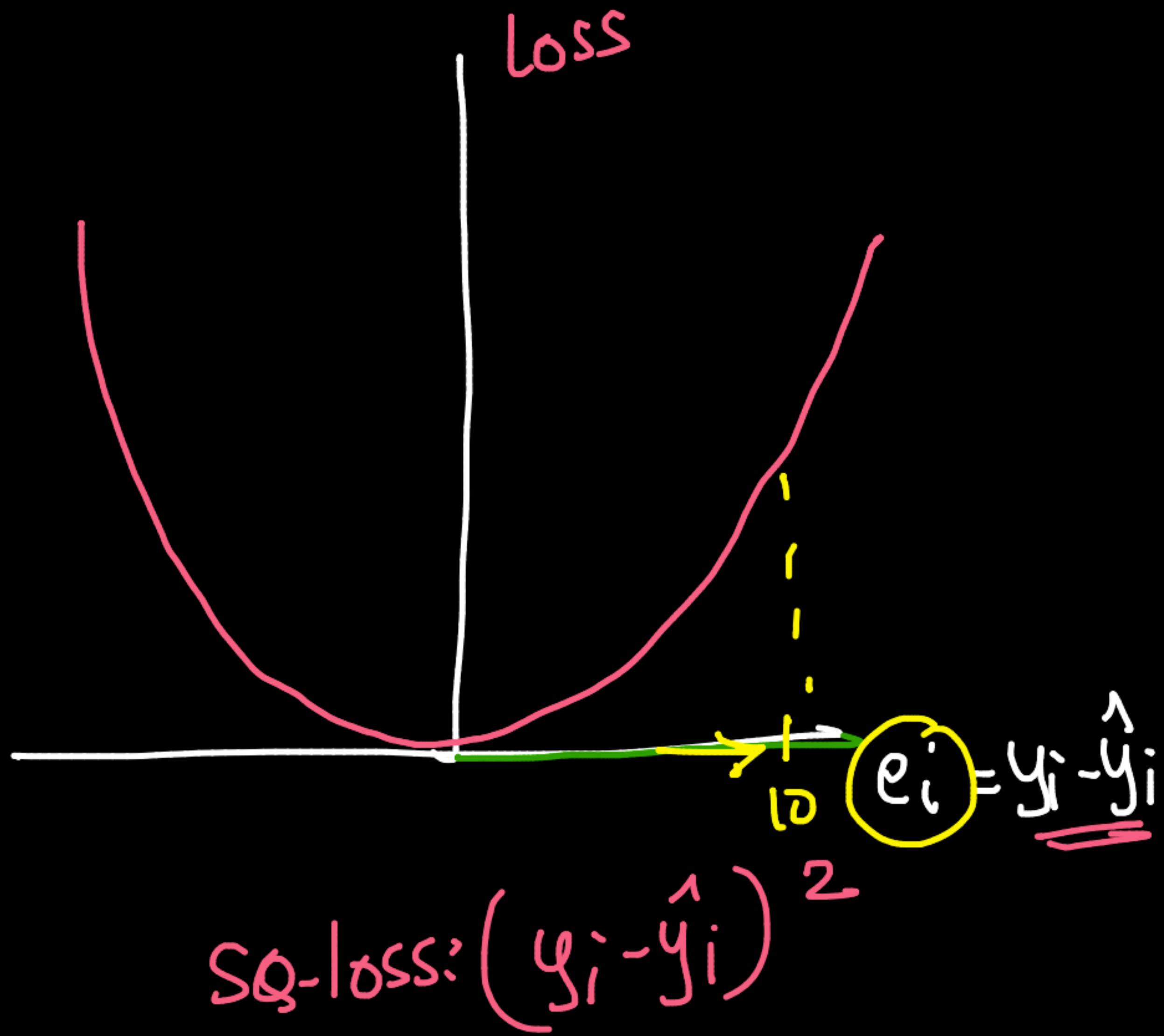
$$\checkmark f_3(x_i) = \circled{1.9} \rightarrow 0.1$$

~~Toy-ex:~~



GBDT : regression

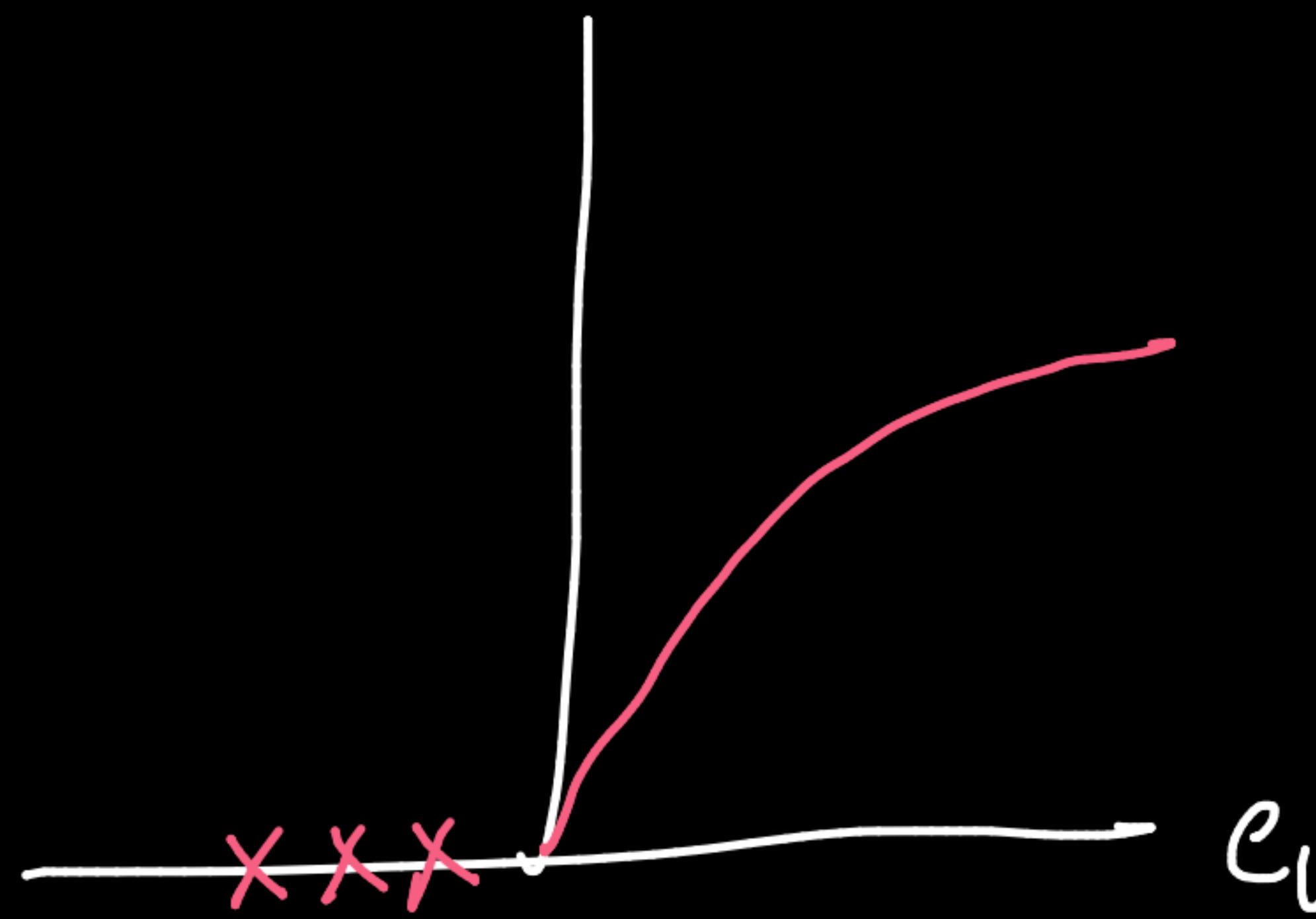




$$\text{Loss} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

x_{100}, y_{100} is an outlier
 $\hat{y}_{100} - y_{100}$: high
error is concentrated
because of the \hat{y}_{100}^2

SQRT



does not work

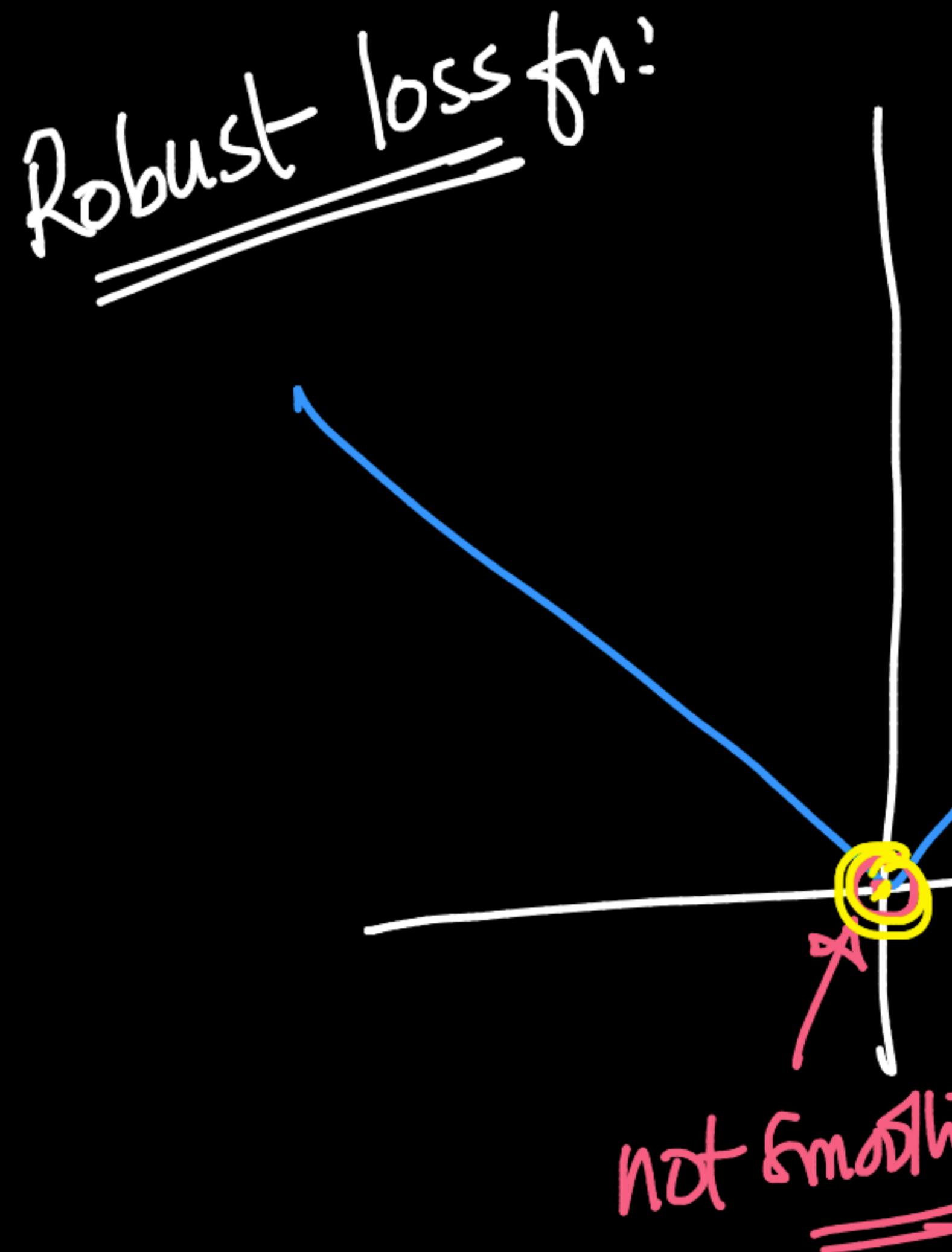


RMSE

better than
Squared loss

$$\sqrt{\frac{1}{n} (e_1^2 + e_2^2 + \dots + e_n^2)}$$

$$e_i = y_i - \hat{y}_i$$



$\text{abs}(e_i)$

$$\min \sum_{i=1}^n f(e_i)$$

$$\text{abs}(x) = \begin{cases} x & \text{if } x \geq 0 \\ -x & \text{if } x < 0 \end{cases}$$

hack

$$\frac{d \text{abs}(x)}{dx} = \begin{cases} +1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$

Math:
NJ

Tools

[What links here](#)[Related changes](#)[Special pages](#)[Permanent link](#)[Page information](#)[Cite this page](#)[Wikidata item](#)[Print/export](#)[Download as PDF](#)[Printable version](#)[Languages](#)[فارسی](#)[Italiano](#)[日本語](#)[Русский](#)[Edit links](#)

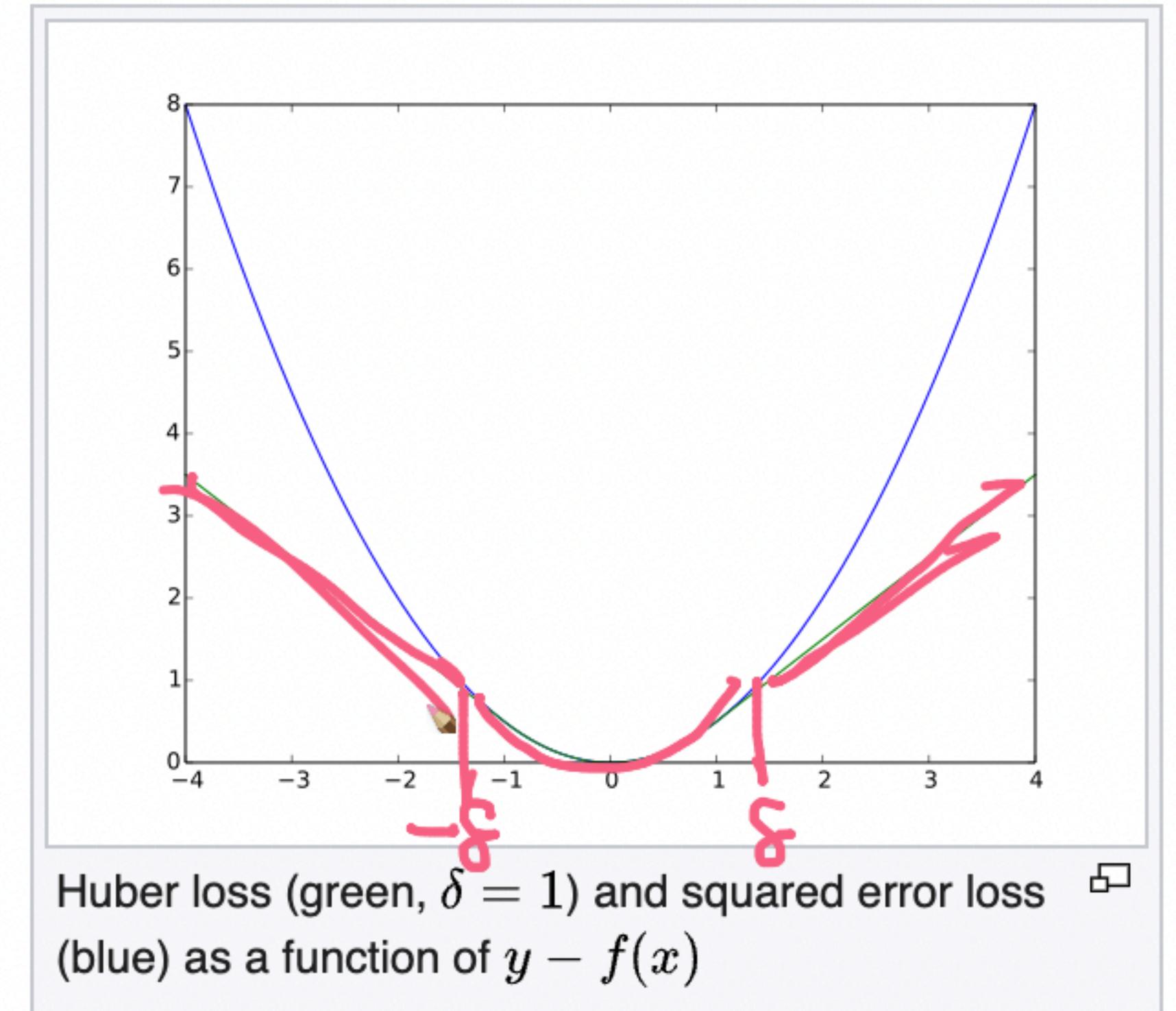
Definition [\[edit \]](#)

The Huber loss function describes the penalty incurred by an [estimation procedure](#). [f. Huber \(1964\)](#) defines the loss function piecewise by^[1]

$$L_\delta(a) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \leq \delta, \\ \delta \cdot \left(|a| - \frac{1}{2}\delta\right), & \text{otherwise.} \end{cases}$$

This function is quadratic for small values of a , and linear for large values, with equal values and slopes of the different sections at the two points where $|a| = \delta$. The variable a often refers to the residuals, that is to the difference between the observed and predicted values $a = y - f(x)$, so the former can be expanded to^[2]

$$L_\delta(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{for } |y - f(x)| \leq \delta, \\ \delta \cdot \left(|y - f(x)| - \frac{1}{2}\delta\right), & \text{otherwise.} \end{cases}$$



Motivation [\[edit \]](#)

Tools

[What links here](#)[Related changes](#)[Special pages](#)[Permanent link](#)[Page information](#)[Cite this page](#)[Wikidata item](#)[Print/export](#)[Download as PDF](#)[Printable version](#)[Languages](#)[فارسی](#)[Italiano](#)[日本語](#)[Русский](#)[Edit links](#)

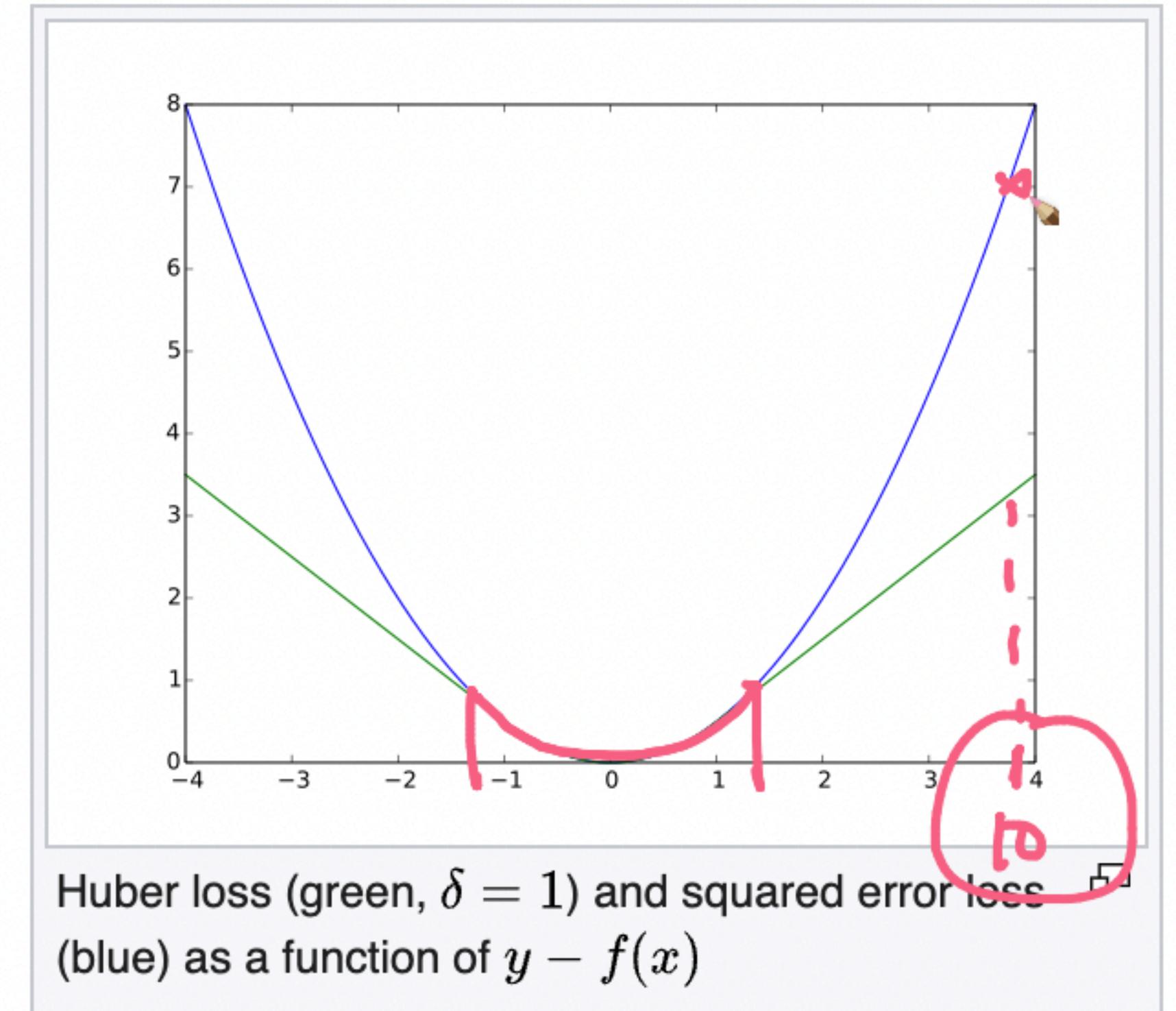
Definition [\[edit \]](#)

The Huber loss function describes the penalty incurred by an [estimation procedure](#). [f. Huber \(1964\)](#) defines the loss function piecewise by^[1]

$$L_\delta(a) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \leq \delta, \\ \delta \cdot \left(|a| - \frac{1}{2}\delta\right), & \text{otherwise.} \end{cases}$$

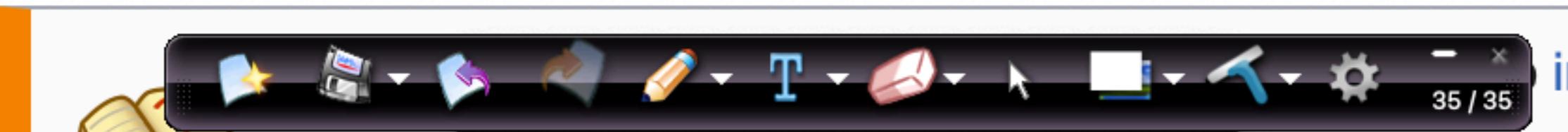
This function is quadratic for small values of a , and linear for large values, with equal values and slopes of the different sections at the two points where $|a| = \delta$. The variable a often refers to the residuals, that is to the difference between the observed and predicted values $a = y - f(x)$, so the former can be expanded to^[2]

$$L_\delta(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{for } |y - f(x)| \leq \delta, \\ \delta \cdot \left(|y - f(x)| - \frac{1}{2}\delta\right), & \text{otherwise.} \end{cases}$$



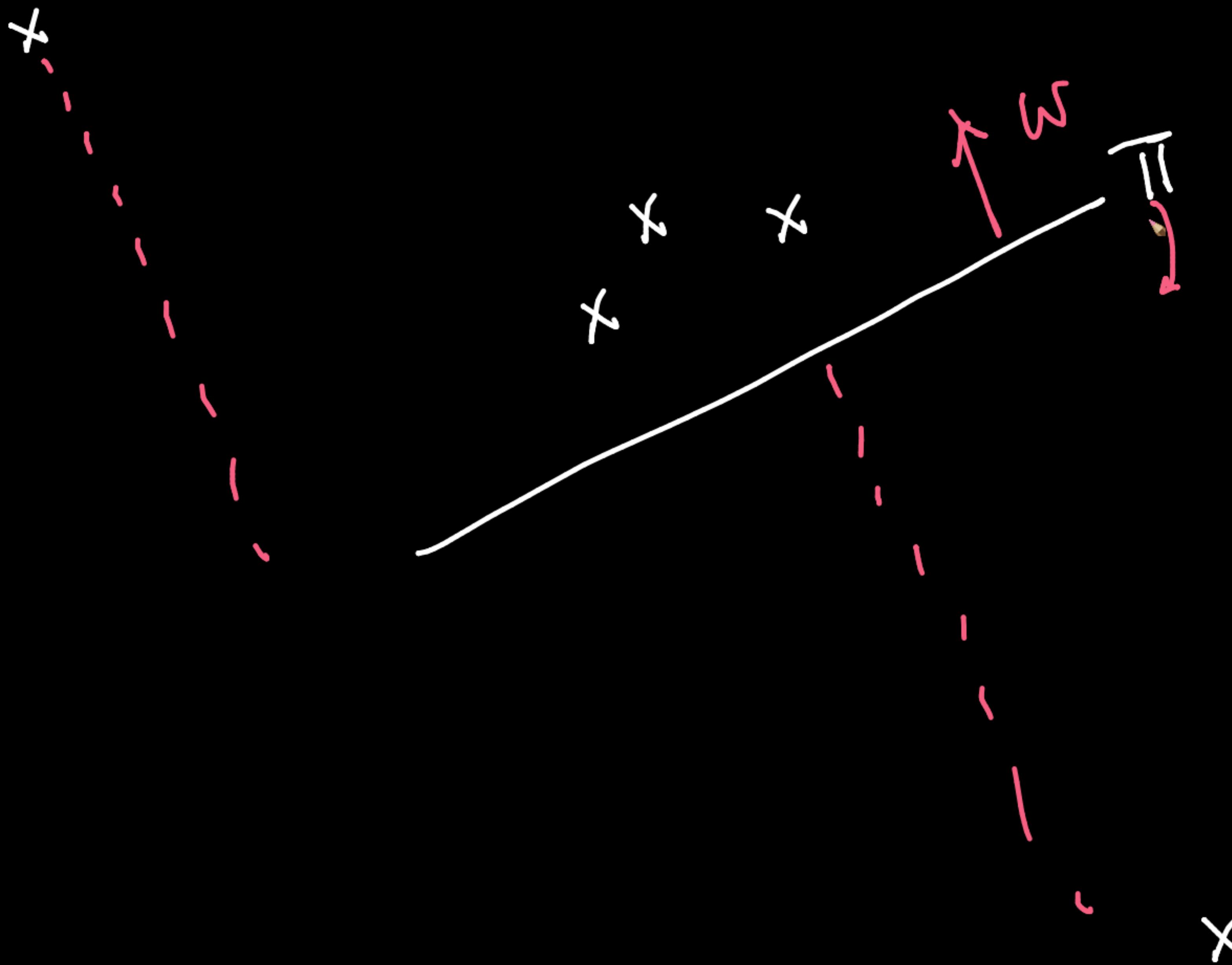
Huber loss (green, $\delta = 1$) and squared error loss (blue) as a function of $y - f(x)$

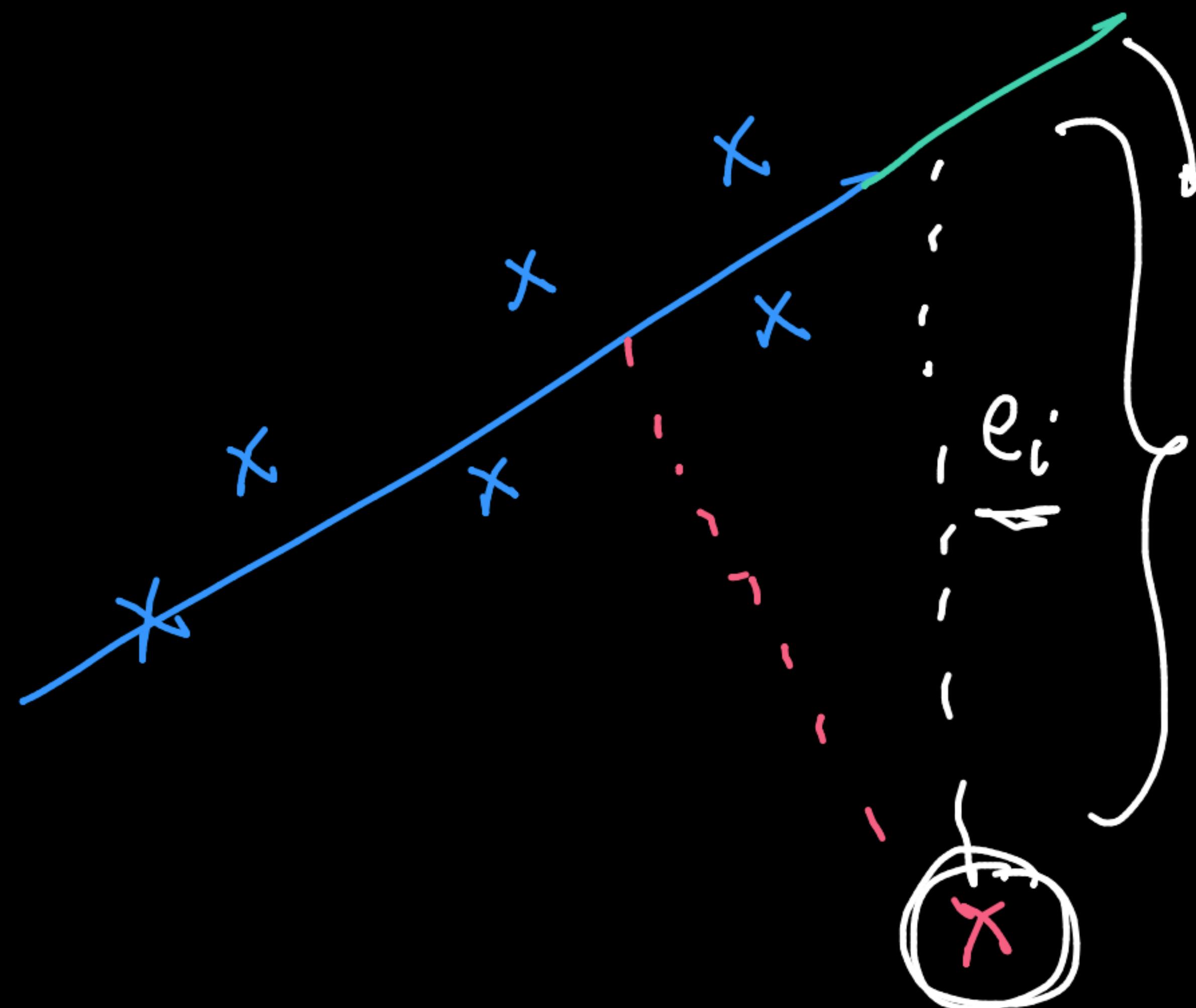
Motivation [\[edit \]](#)



improve this section by

adding citations to reliable sources. Unourced material may be challenged and removed.





huber-loss

A hand-drawn graph of the Huber loss function. It is a piecewise function with a linear slope for small values and a parabolic shape for larger values, transitioning at zero.

ESLII.pdf XgBoost.ipynb - Colaboratory | Parameters – LightGBM 3.3.2. | Distributed Learning Guide – L | Huber loss - Wikipedia

hastie.su.domains/Papers/ESLII.pdf

379 / 764 | - 200% + | ☰ ⌂

ESLII.pdf

360 10. Boosting and Additive Trees

TABLE 10.2. *Gradients for commonly used loss functions.*

Setting	Loss Function	$-\partial L(y_i, f(x_i))/\partial f(x_i)$
Regression	$\frac{1}{2}[y_i - f(x_i)]^2$	$y_i - f(x_i)$
Regression	$ y_i - f(x_i) $	$\text{sign}[y_i - f(x_i)]$
Regression	Huber	$y_i - f(x_i)$ for $ y_i - f(x_i) \leq \delta_m$ $\delta_m \text{sign}[y_i - f(x_i)]$ for $ y_i - f(x_i) > \delta_m$ where $\delta_m = \alpha\text{th-quantile}\{ y_i - f(x_i) \}$
Classification	Deviance	$k\text{th component: } I(y_i = \mathcal{G}_k) - p_k(x_i)$

boosting procedure, and top-down decision tree induction, are themselves approximation procedures. After constructing the tree (10.37), the corresponding constants in each region are given by (10.30).

Table 10.2 summarizes the gradients for commonly used loss functions. For squared error loss, the negative gradient is just the ordinary residual $-g_{im} = y_i - f_{m-1}(x_i)$, so that (10.37) on its own is equivalent standard least squares boosting. With absolute error loss, the negative gradient is the *sign* of the residual, so at each iteration (10.37) fits the tree to the sign of the current residuals by least squares. For Huber M-regression, the

table).

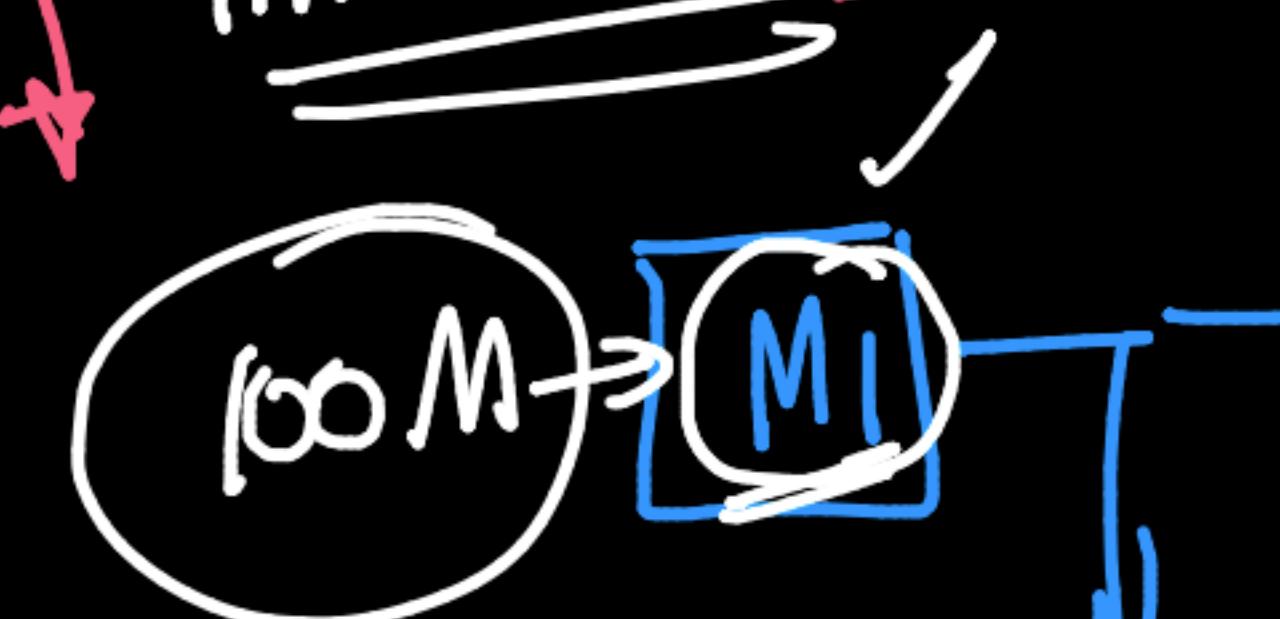
ance (10.22),

Cascading

Fraud-detection

TRMS
credit

imbalanced



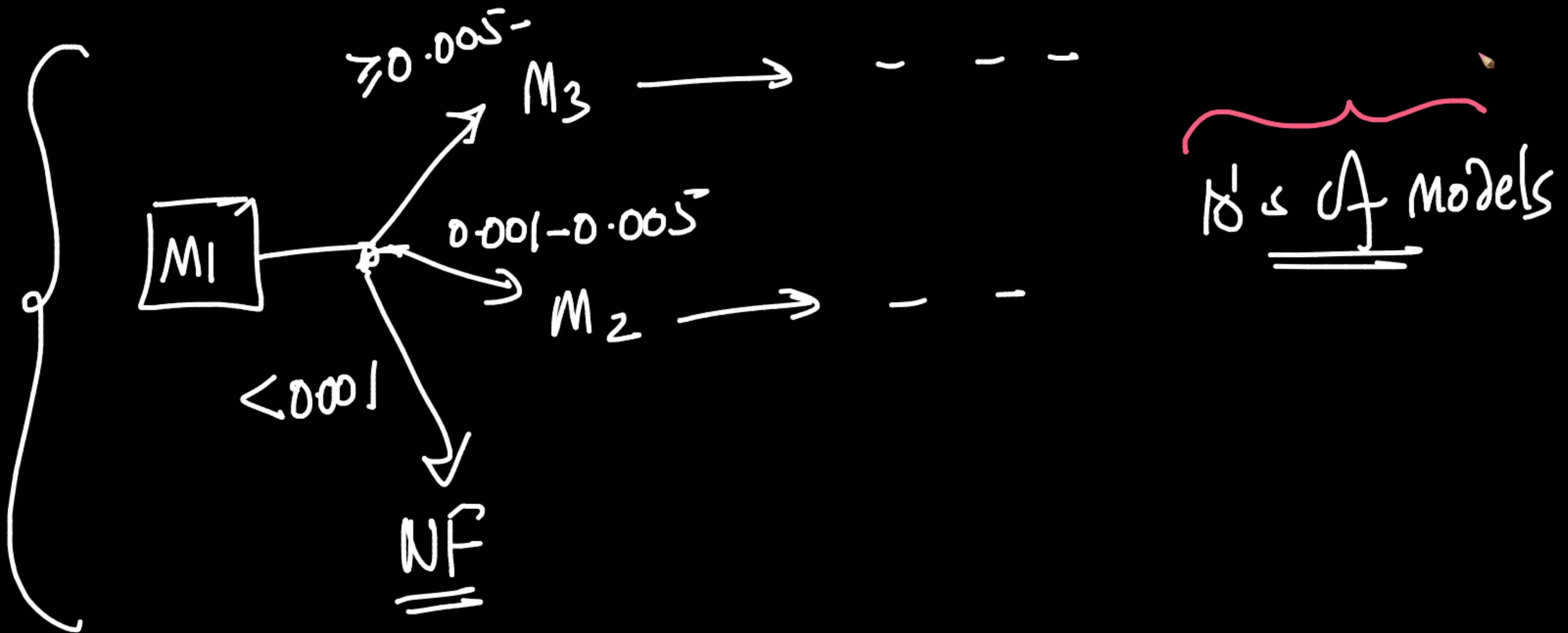
Stacking

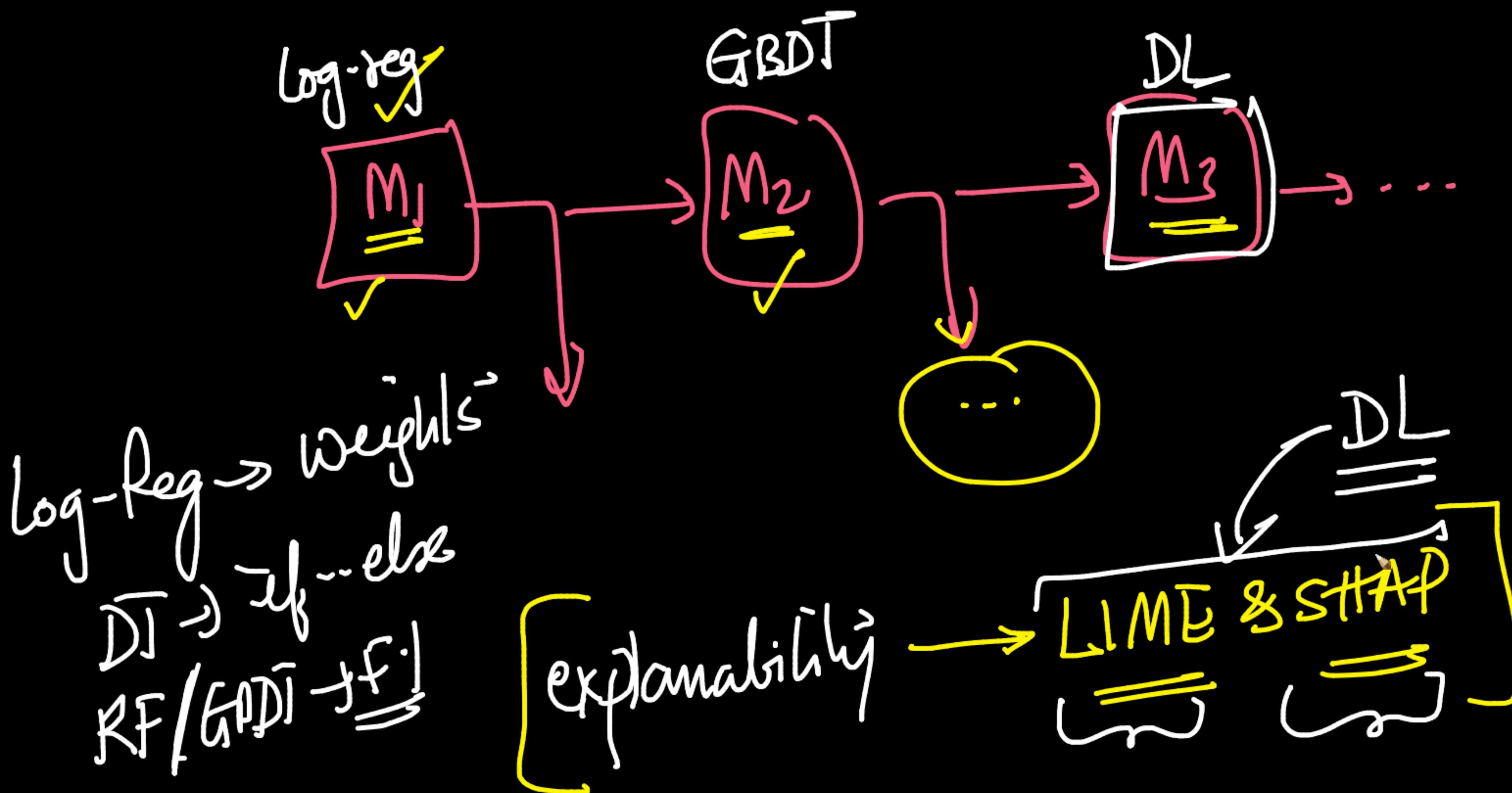
Kaggle - comp

AUC: →

1st : 0.8234
2nd : 0.8231

DL learning





Kaggle vs Real-world

→ Discussion froms & Winners Sols...

→ One-metric → Squeeze out value

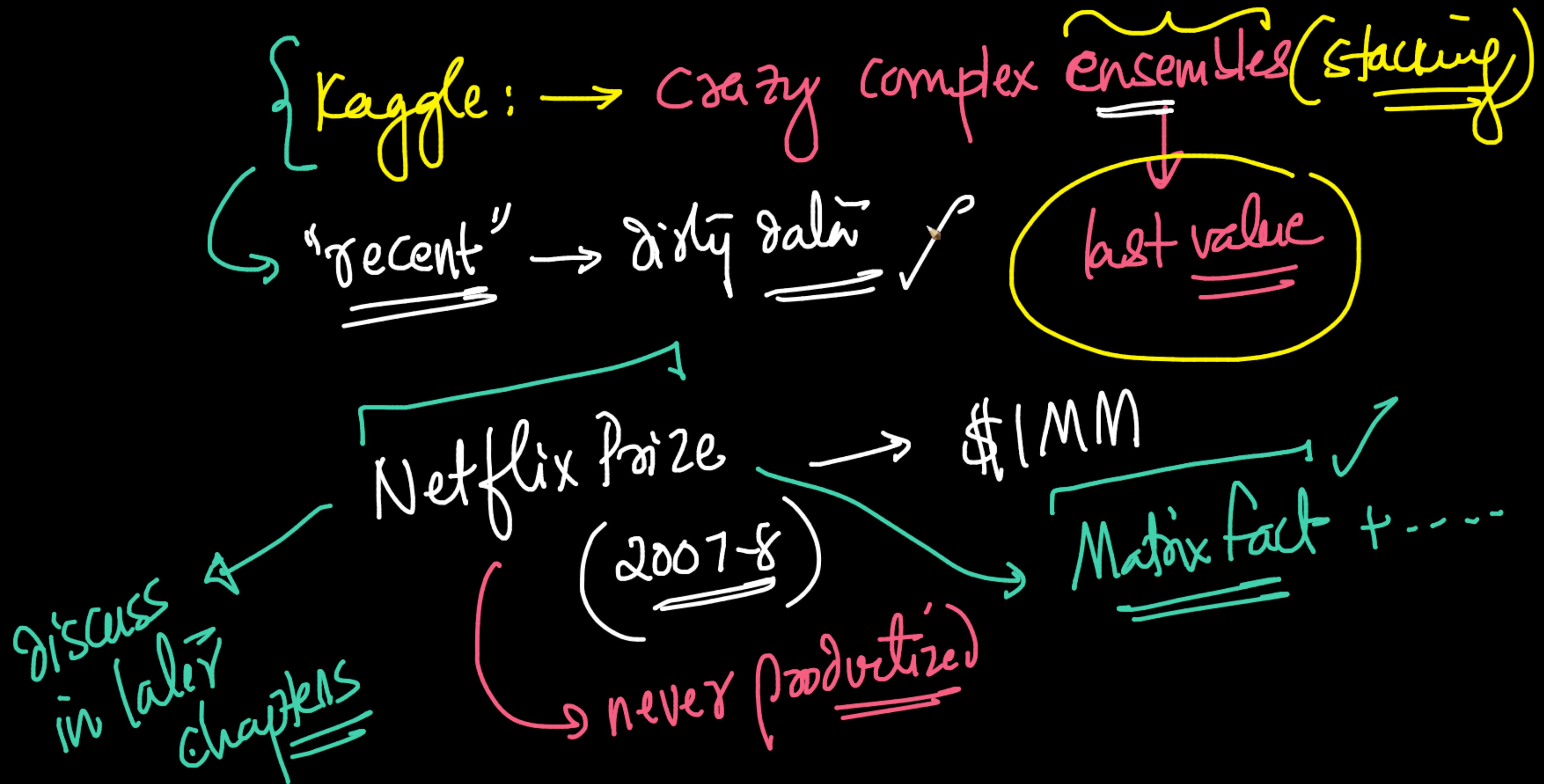
→ multiple things

→ Interpretability

→ Speed / latency

metrics { F1 &
bit-metrics







↳ 1990's

USSR (60's ≤ 80's)

SVM
Kernel Methods

↳ ... theoretical ...

↳ less frequently used now a days

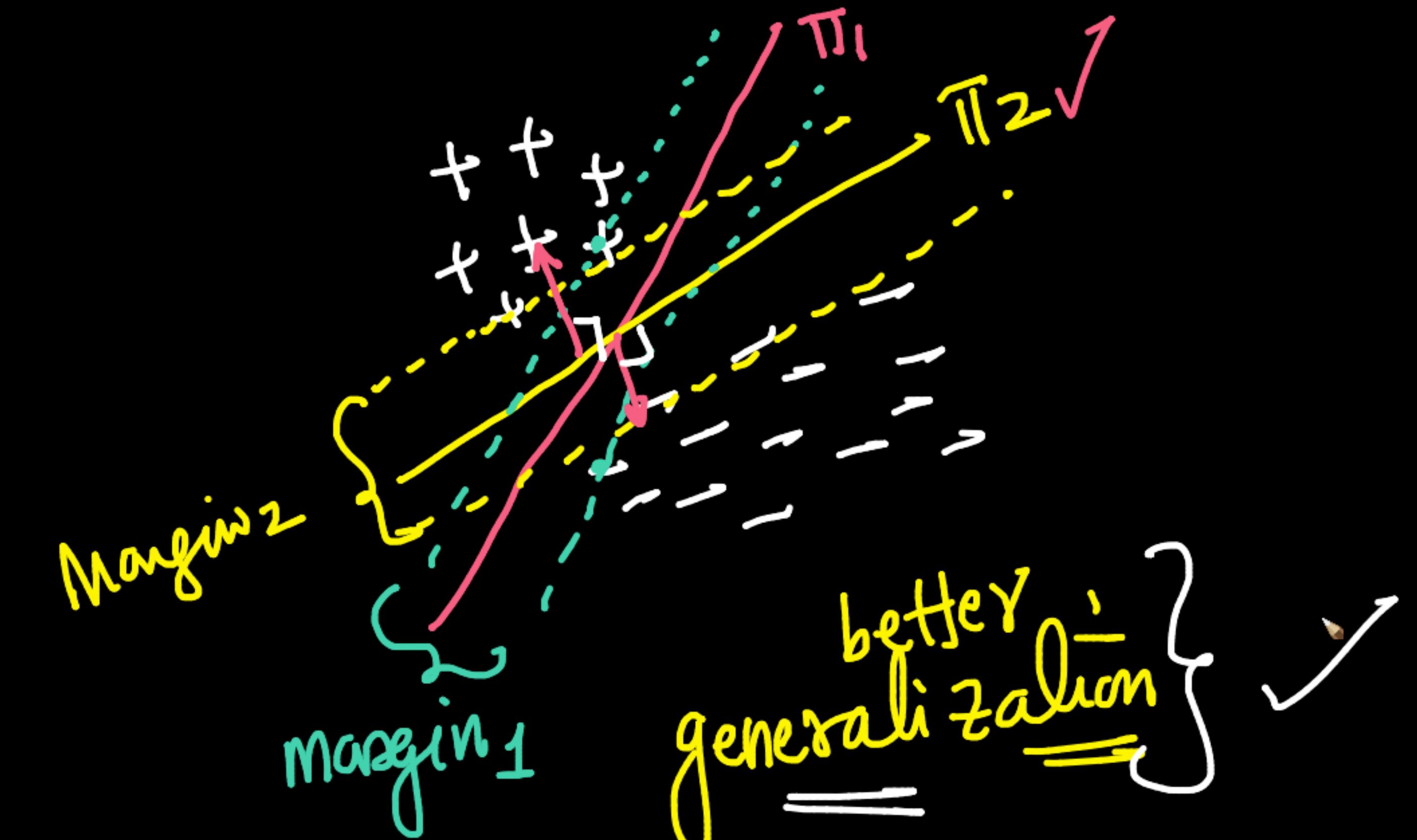
→ one of the more challenging algo

Math

Geom:

{ pick π that
results in
largest margin

{ margin-Maximising
Classifer

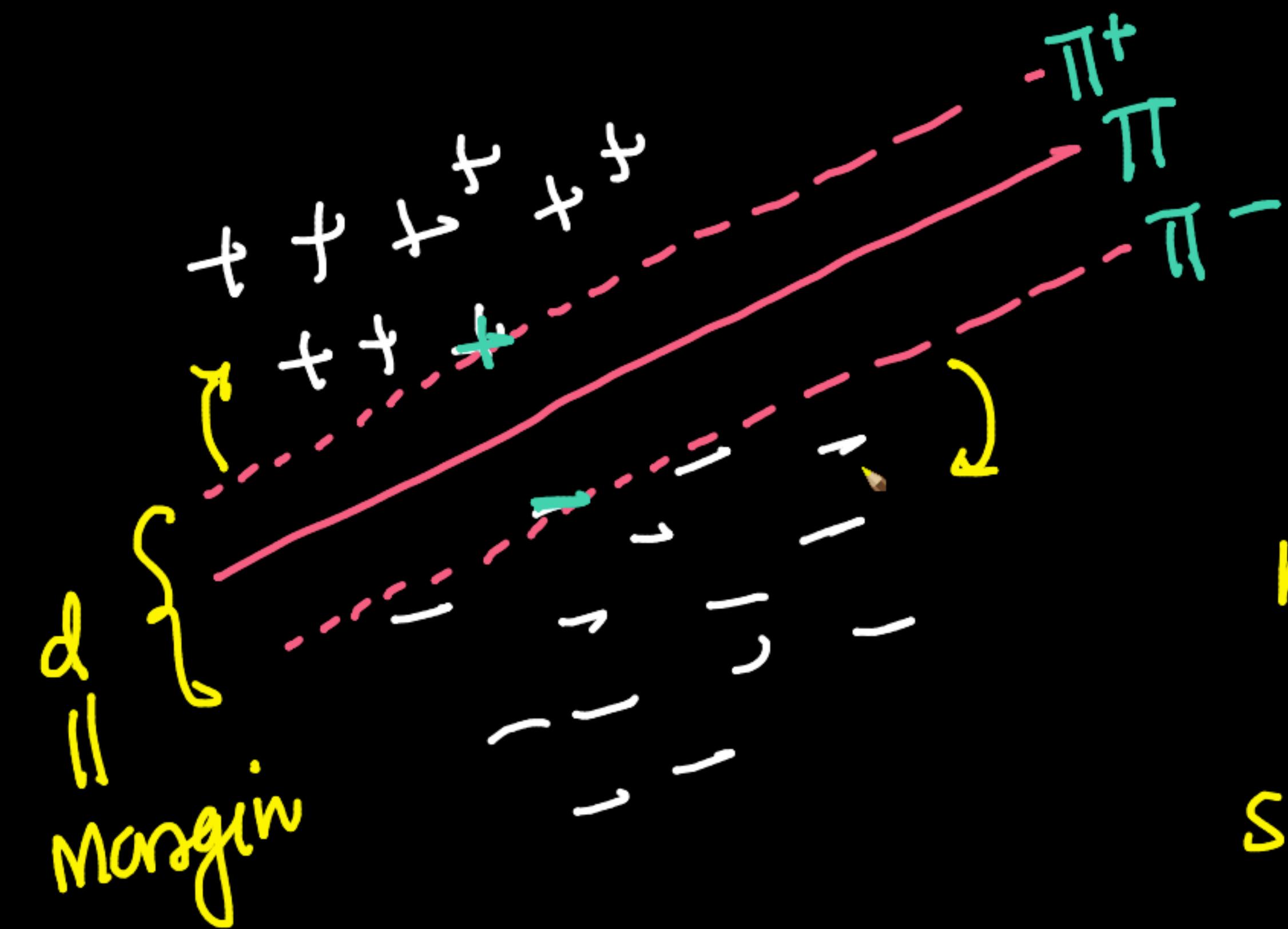


better
generalization

π : Separating hyperplane

$\pi^+ \parallel \pi$

$\pi^- \parallel \pi$



Max. Margin
 $s.t.$

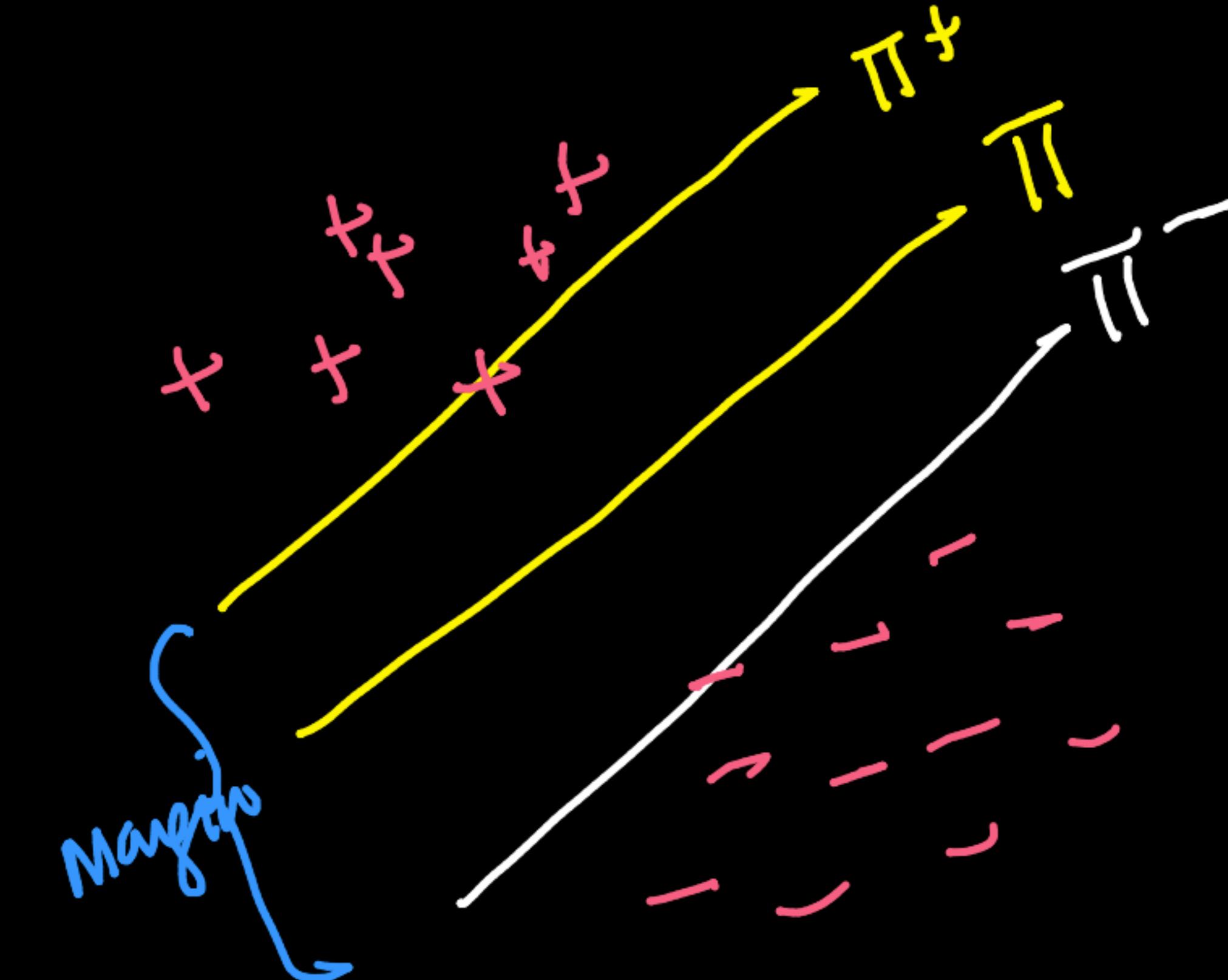
(Let)

$$\Pi: \underline{w^T x + b = 0}$$

$$\Pi^+: w^T x + b = 1$$

$$\Pi^-: w^T x + b = -1$$

$$\|w\| \neq 1$$



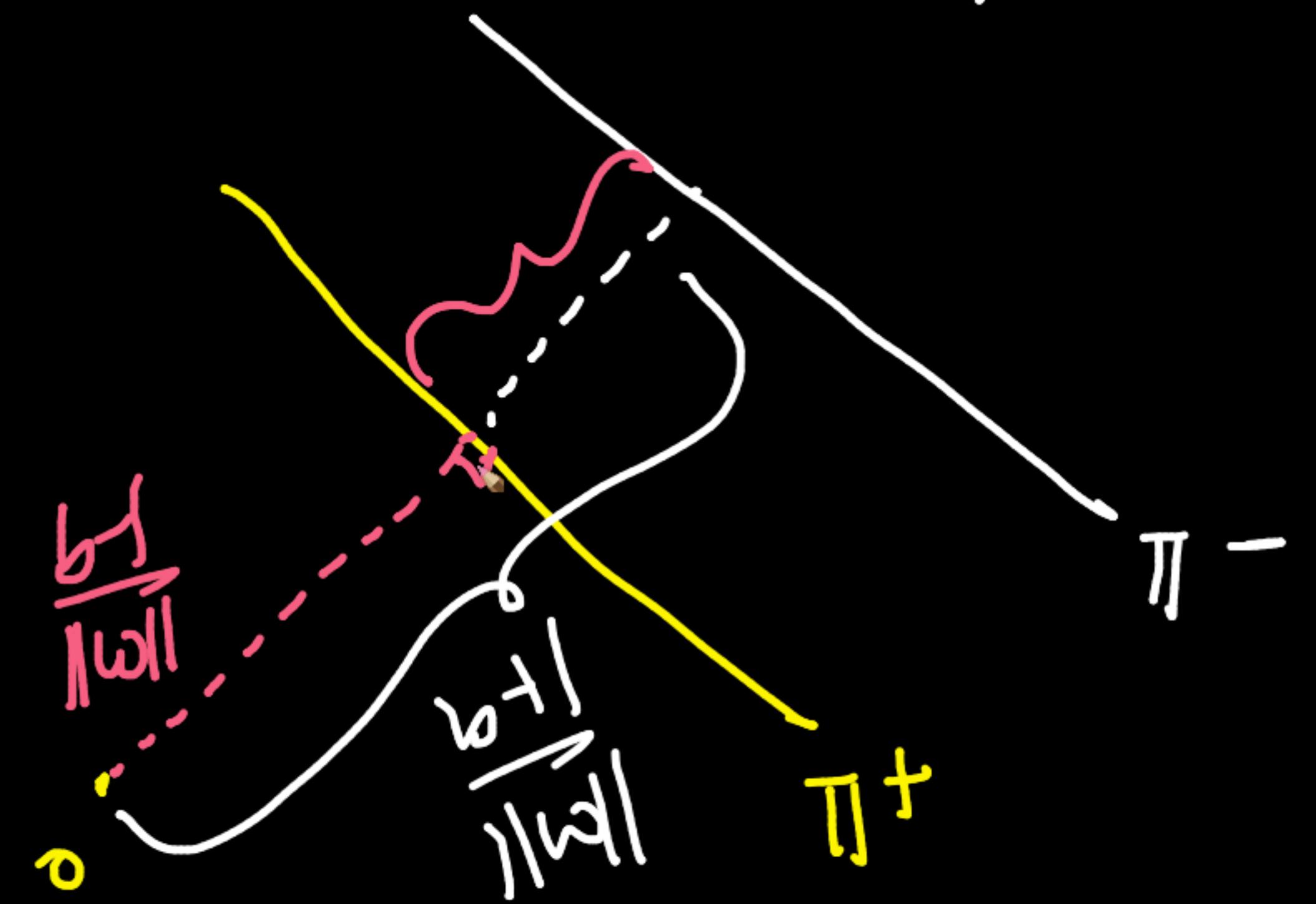
$$\text{Margin} = \frac{2}{\|w\|}$$

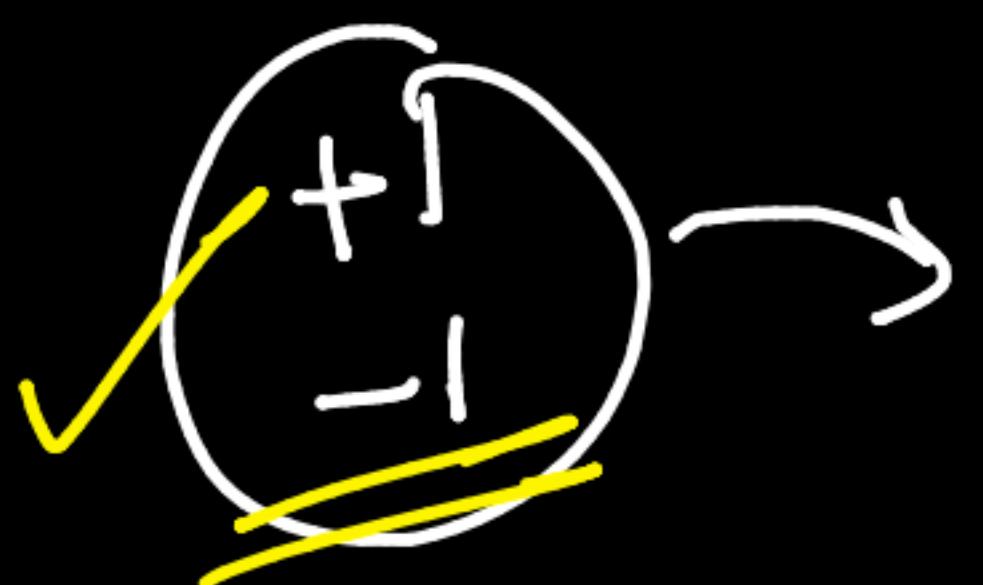
Simple Case: Linearily Separable Data

$$\pi^+ : \tilde{w}^\top x + \tilde{b} \stackrel{\sim}{=} 0$$

$$\pi^- : \tilde{w}^\top x + \tilde{b} \stackrel{\sim}{=} 0$$

$$d(\pi^+, \pi^-) = \frac{2}{\|\tilde{w}\|}$$

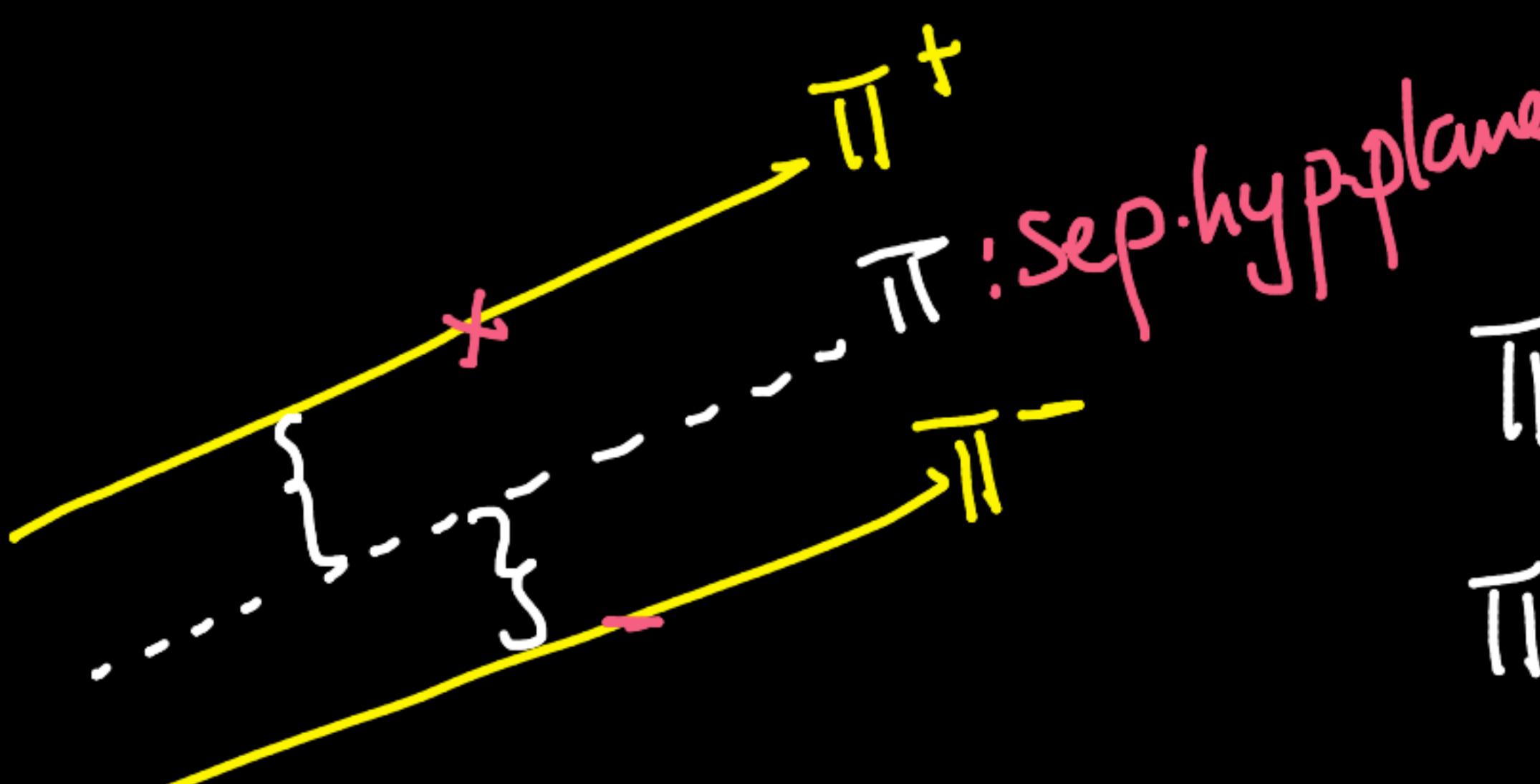
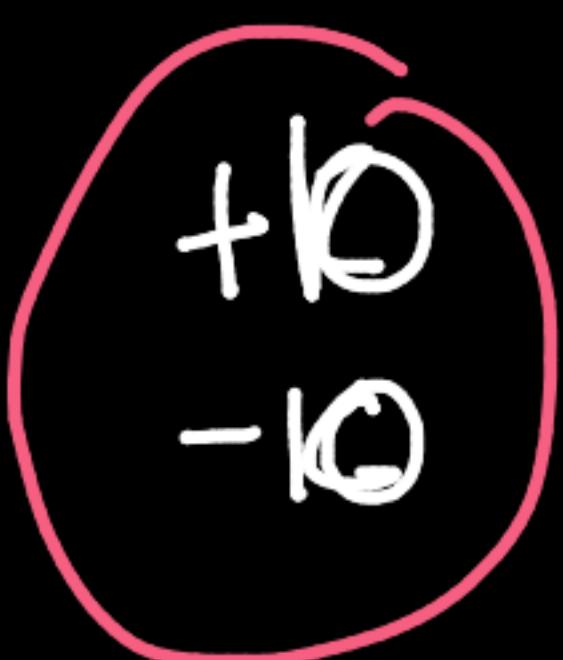




$$\text{Max}_{w,b} \frac{2}{\|w\|}$$



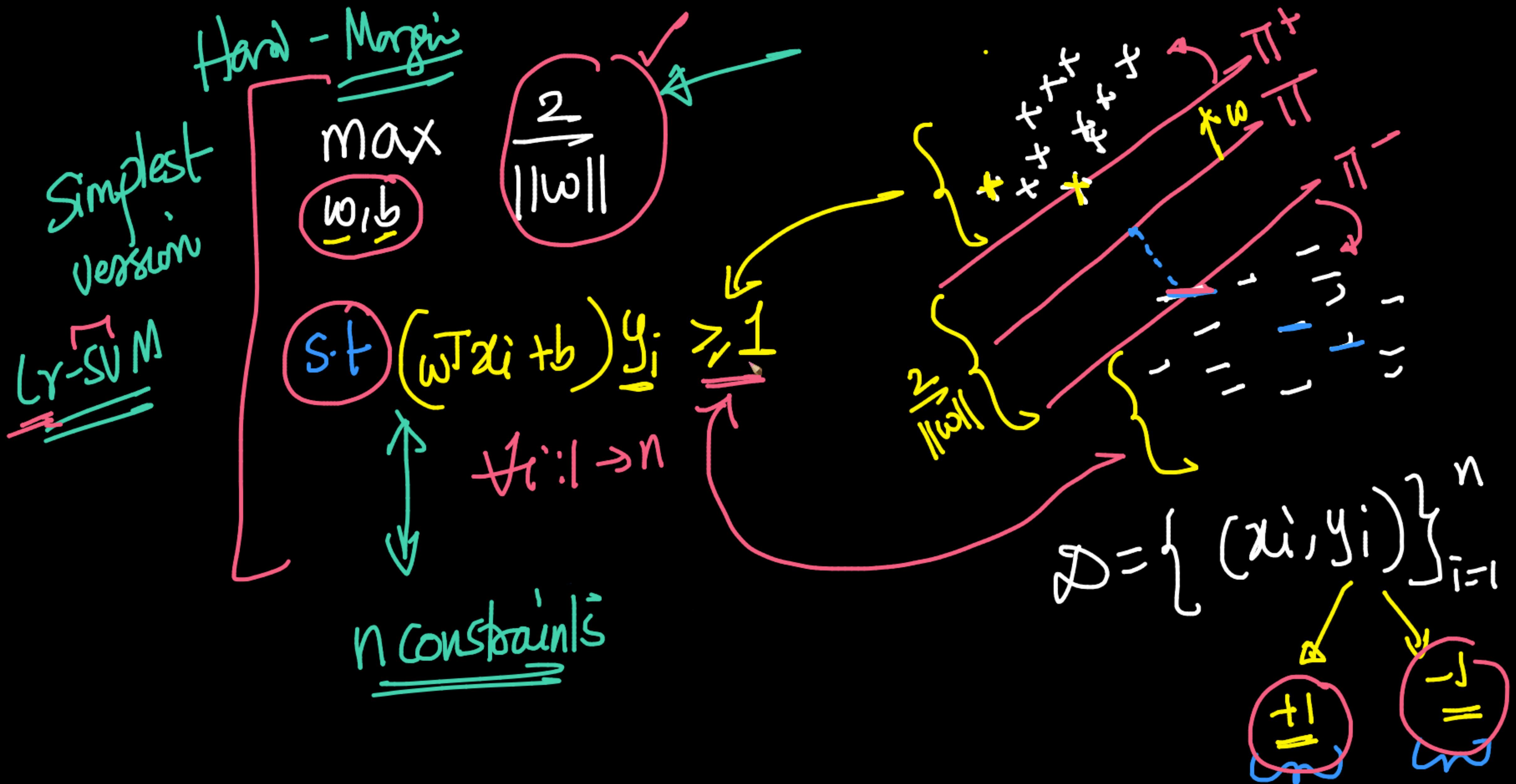
$$\text{Max}_{w,b} \frac{20}{\|w\|}$$



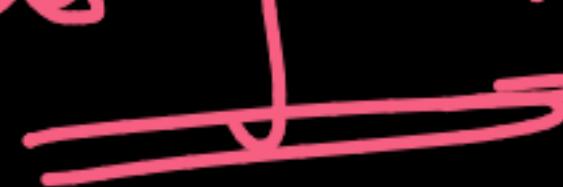
$$\pi^+: w^T x + b = +10$$

$$\pi^-: w^T x + b = -10$$

$$\text{Margin} = \frac{20}{\|w\|}$$

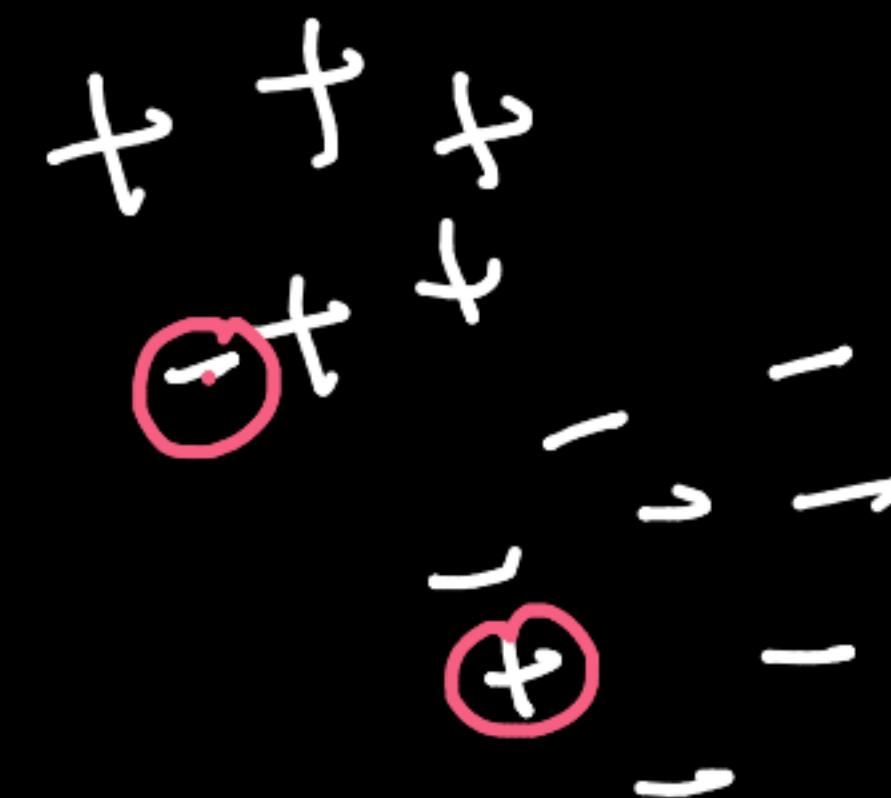


When would this fail: Hard Margin SVM:-

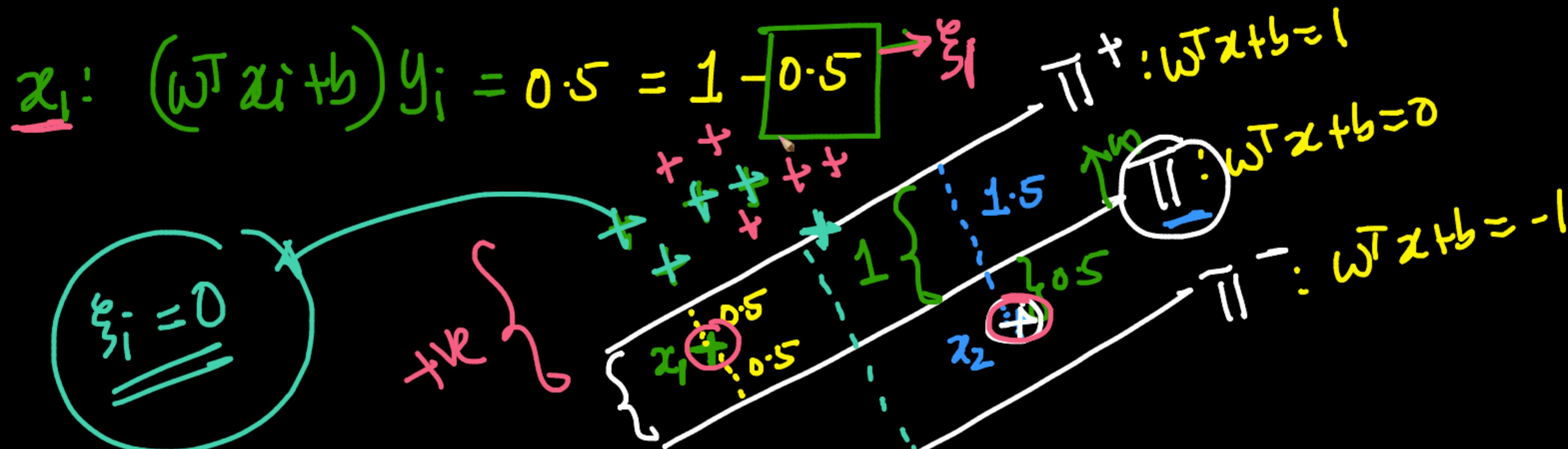


↓ Data is non-linear.

Sep:



$$\pi: \underline{\theta}^T x + \underline{\theta}_0$$



$$\chi_2: (\bar{w}^T \bar{x}_i + b) y_i = -0.5 = 1 - 1.5 \rightarrow \xi_2$$

-ve

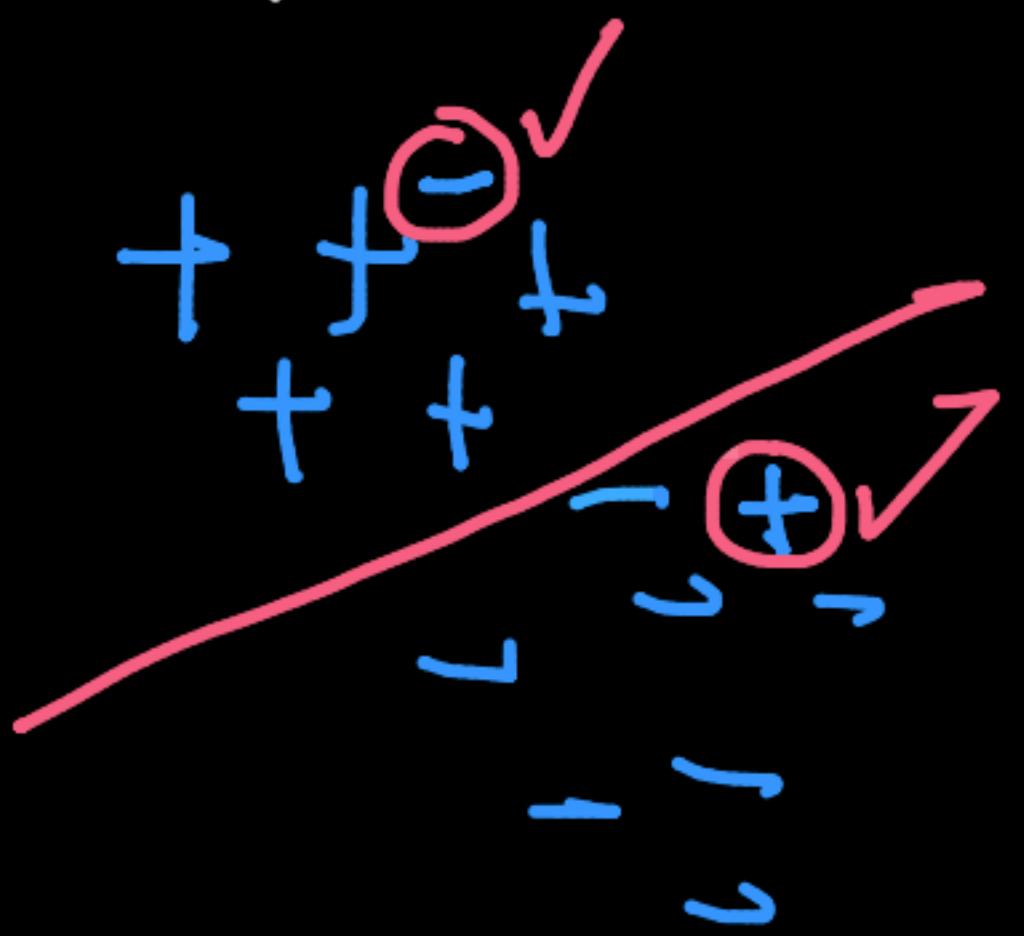
$$\chi_3: (\bar{w}^T \bar{x}_i + b) y_i = 1 - 2.5 \rightarrow \xi_3$$

- $\xi_i = 0$ for all 'correctly' placed pls \rightarrow hard-Margin
- $\xi_i > 0$ for all 'incorrectly' placed pls
- $\hookrightarrow \xi_i < 1 \rightarrow$ it will still be classified correctly by π
- $\xi_i > 1 \rightarrow$ incorrectly classified

Soft-Margin SVM

$$\text{Max } \frac{2}{\|\omega\|} = \text{Min } \frac{\|\omega\|}{z}$$

not perfectly
L₂-separable



almost linear
Separable

Soft Margin

LR-SVM

$$\min_{w,b} \left[\frac{\|w\|}{2} + C \sum_{i=1}^n \xi_i \right]$$

HP

s.t.

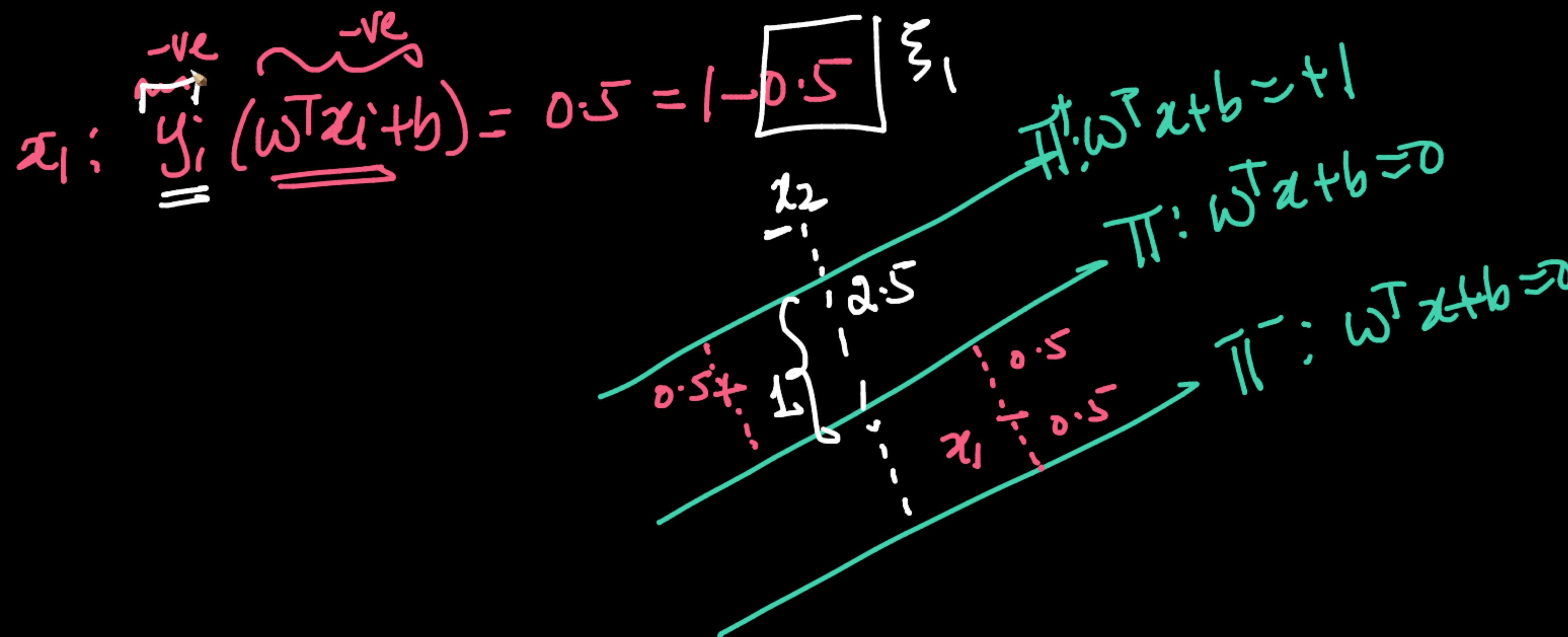
$$y_i (\bar{w}^T x_i + b) \geq 1 - \xi_i$$

$\forall i: 1 \rightarrow n$

$$\xi_i \geq 0 \quad \forall i: 1 \rightarrow n$$

for +ve &
-ve pts

Linear Classifier
[almost LR-Separator]

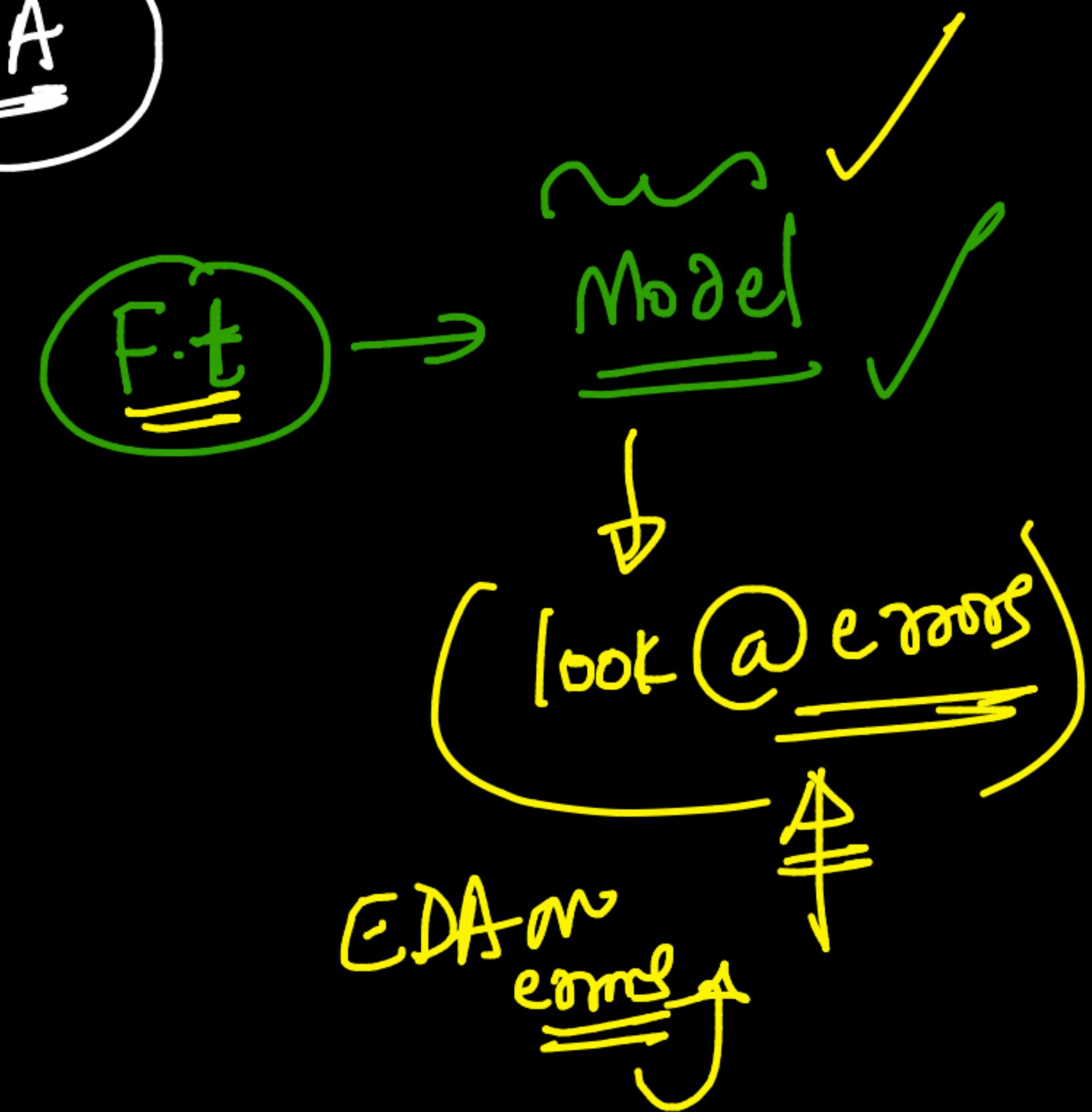


$$x_2: y_2 (\bar{w}^T \bar{x}_2 + b) = 2.5$$
$$= 1 - \boxed{2.5} \xi_2$$

ML - problem

→ Understand: EDA

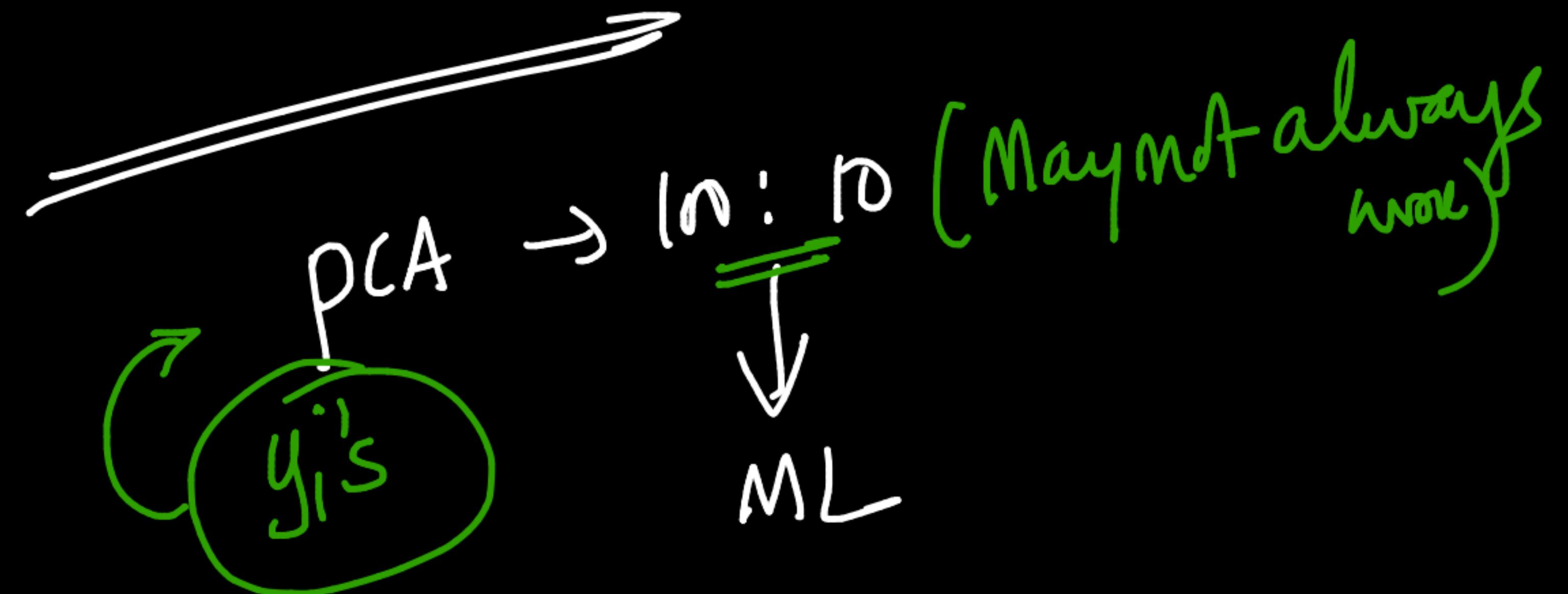
→ Mapping:



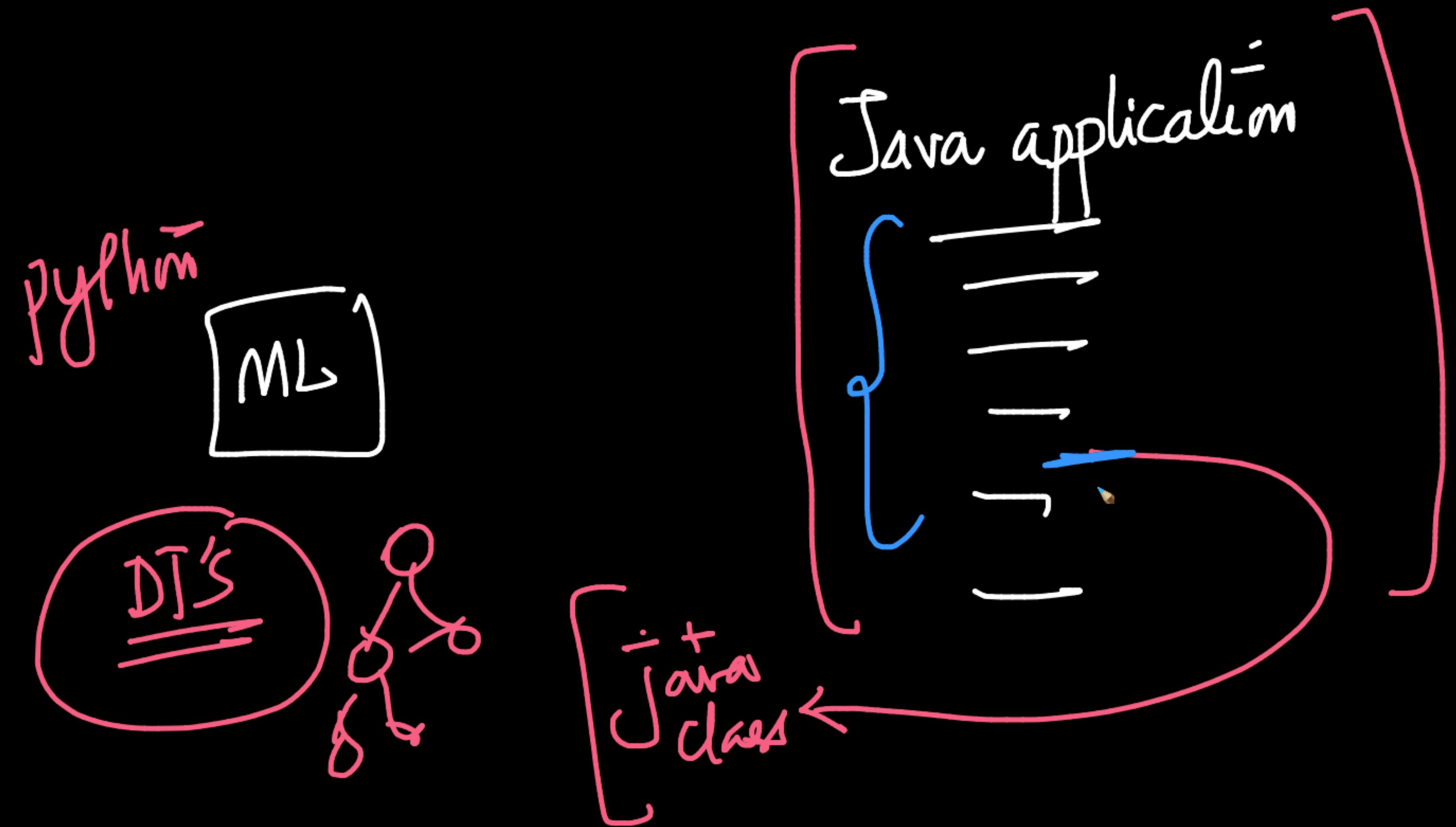
✓ linear Model

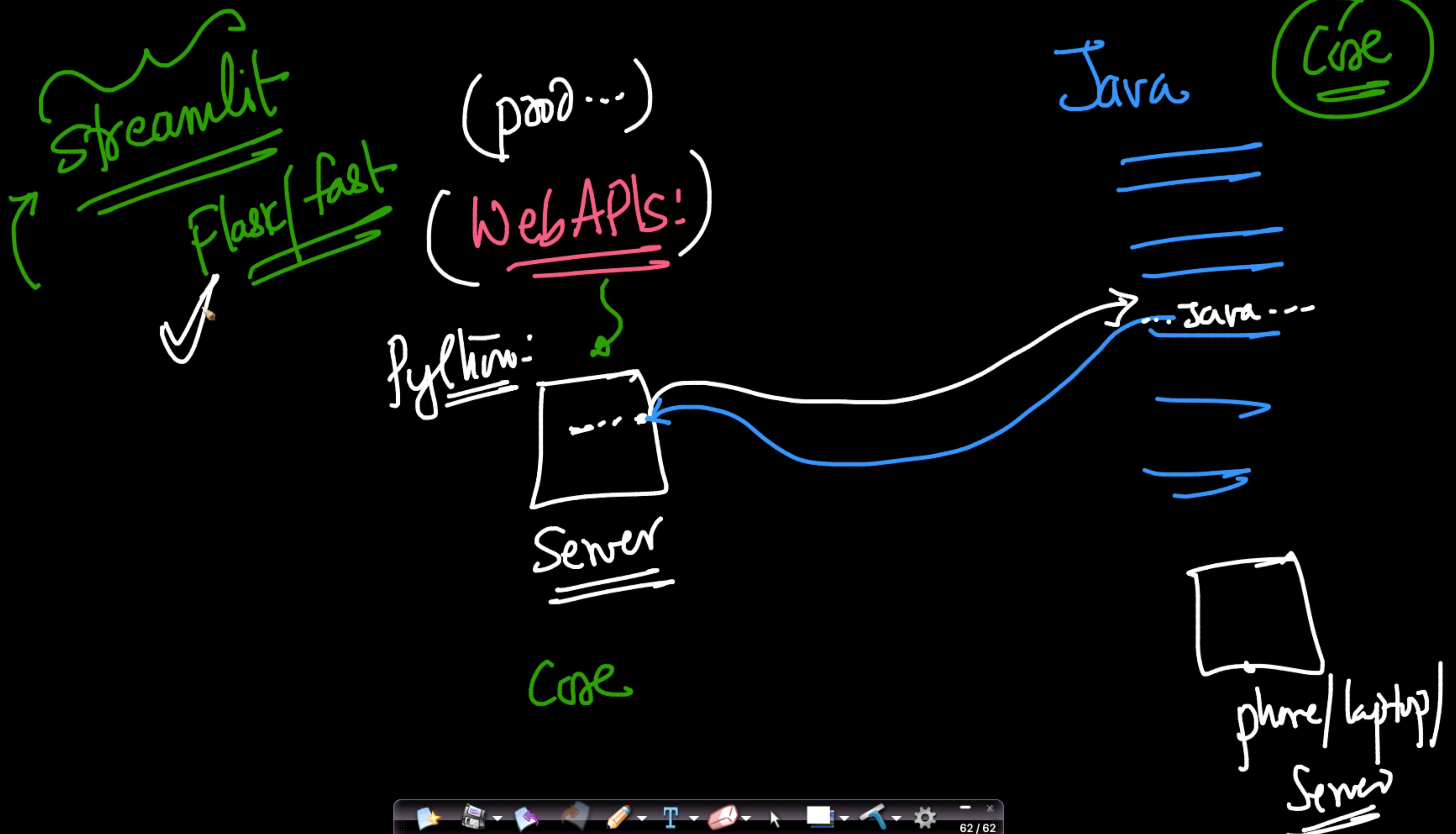
{ $d=10$
 $n=1000$

"regularize" - heavily ✓





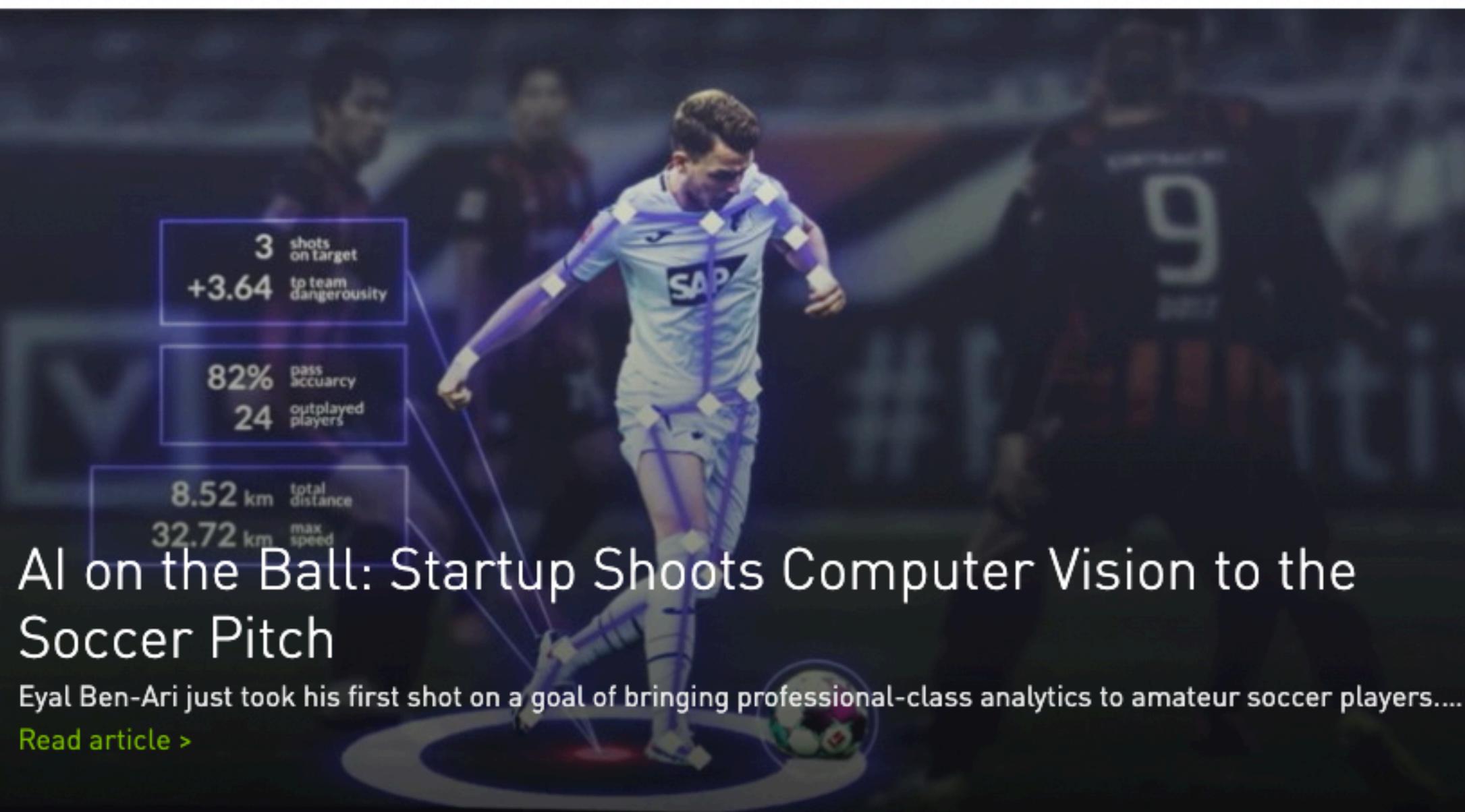




Pyth⁺
C++: Numpy
=
(not a code)

[HOME](#) [AI](#) [DATA CENTER](#) [DRIVING](#) [GAMING](#) [PRO GRAPHICS](#) [AUTONOMOUS MACHINES](#) [HEALTHCARE](#) [STARTUPS](#) [AI PODCAST](#)

DEEP LEARNING



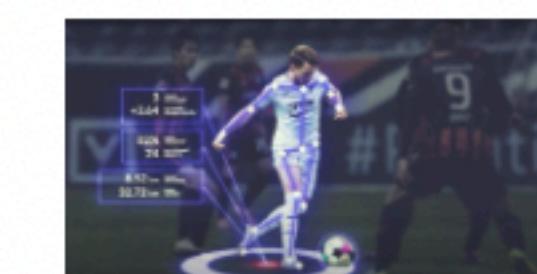
MOST POPULAR



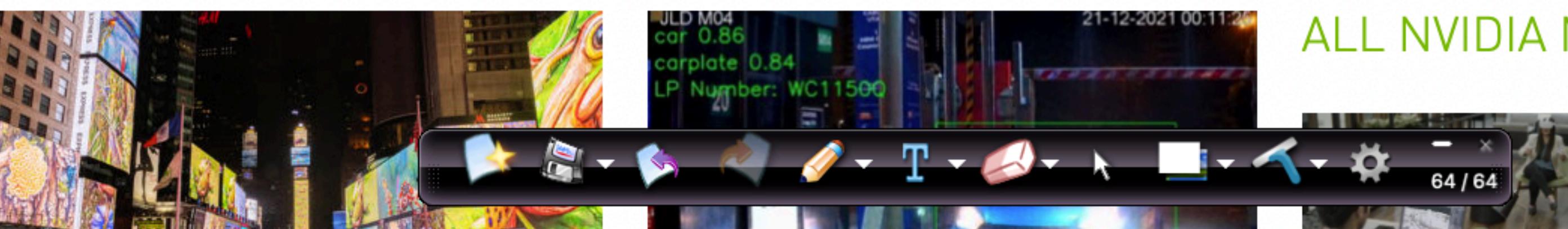
What is Extended Reality?



Mission Made Possible: Real-Time Rendering Helps Studio Create Cinematic Battle Between Characters From 'Diablo Immortal'



AI on the Ball: Startup Shoots Computer Vision to the Soccer Pitch



ALL NVIDIA NEWS

What is Extended Reality?

