# Actual data

2M Rows ⟶ 50 - 100 coloums

- Linear model
- KNN

- DT
- RF
- Ada Boost

(4-5 algo)

## Naive Bayes ⟶

10k Rows and 5k coloums

Small data size
10k and 20k — 50k

Linear models and (KNN)
good accuracy



Linearly separable data.

Logistic reg.

$$ y = \frac{1}{1 + e^{y}} $$

0 to 1

100 %

Safron    Green

① ② ③

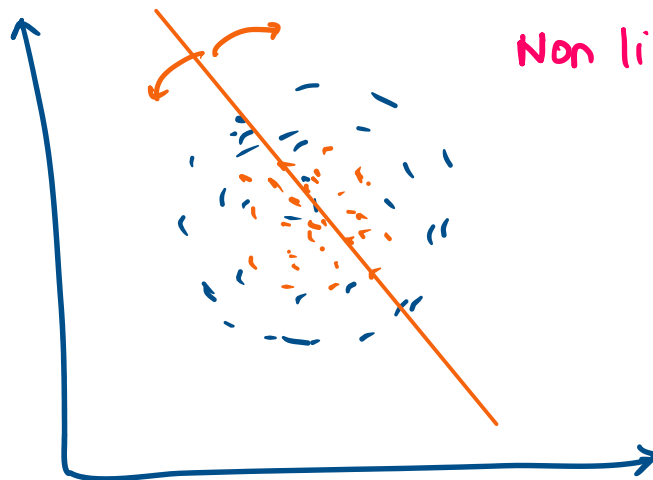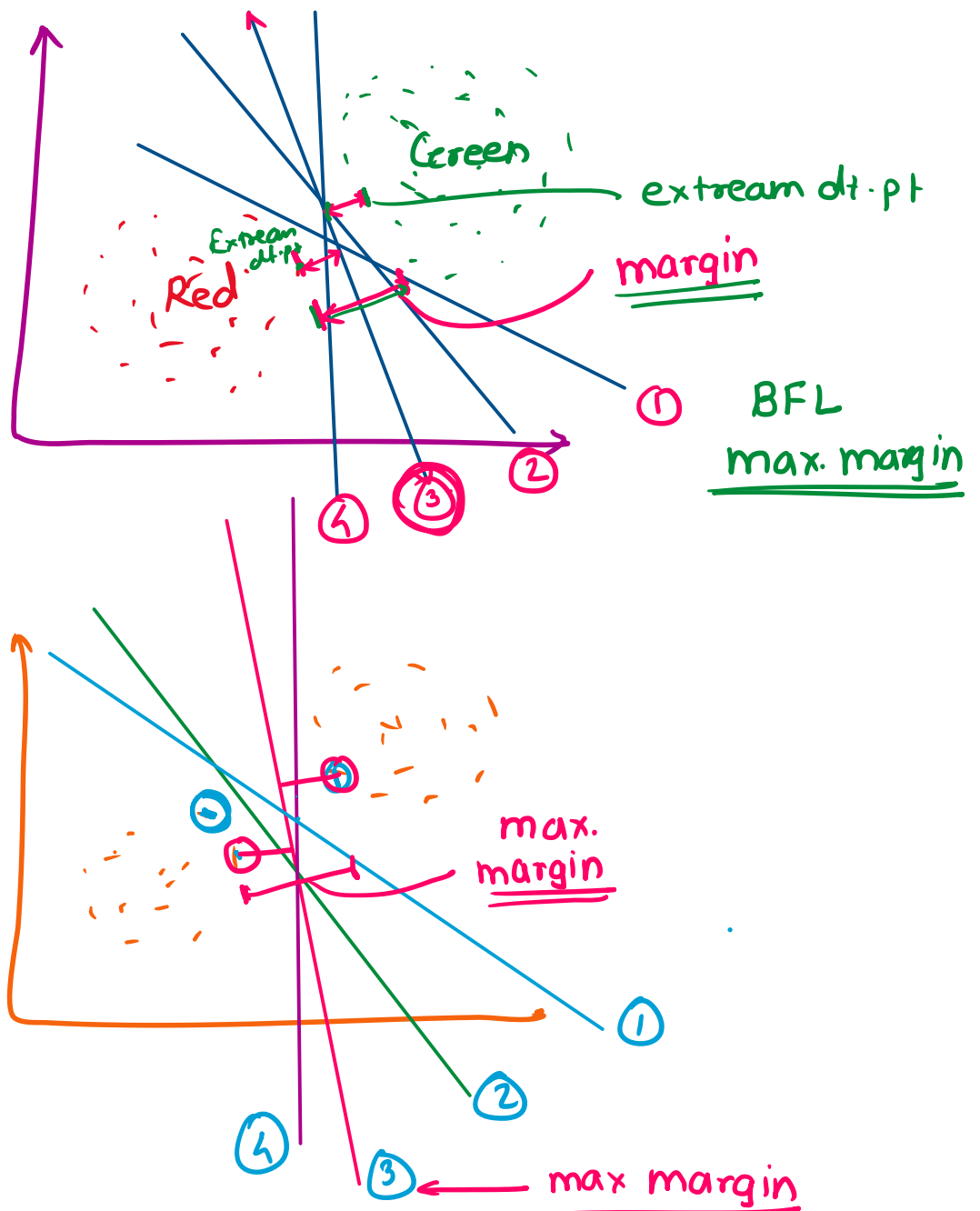Linear line

_Linear model
∨ KNN
_DT
✓ Ensemble



Non linearly separable

SVM will work.

SVM → Support Vector Machine

   handle both linearly separable and
non linearly separable data

① linearly separable data

Green

extream dt.pt

margin

Red

Extream
dt.pt

① BFL
max. margin

②

③

④

max.
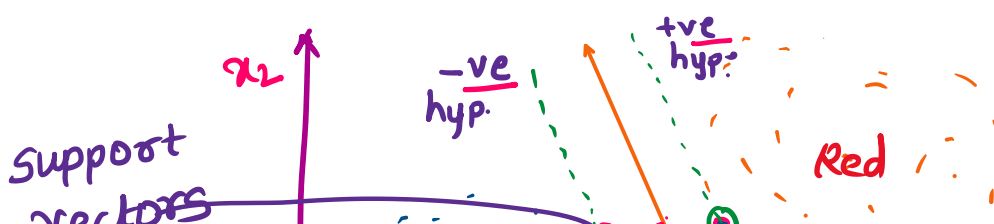margin

①

②

③ ← max margin
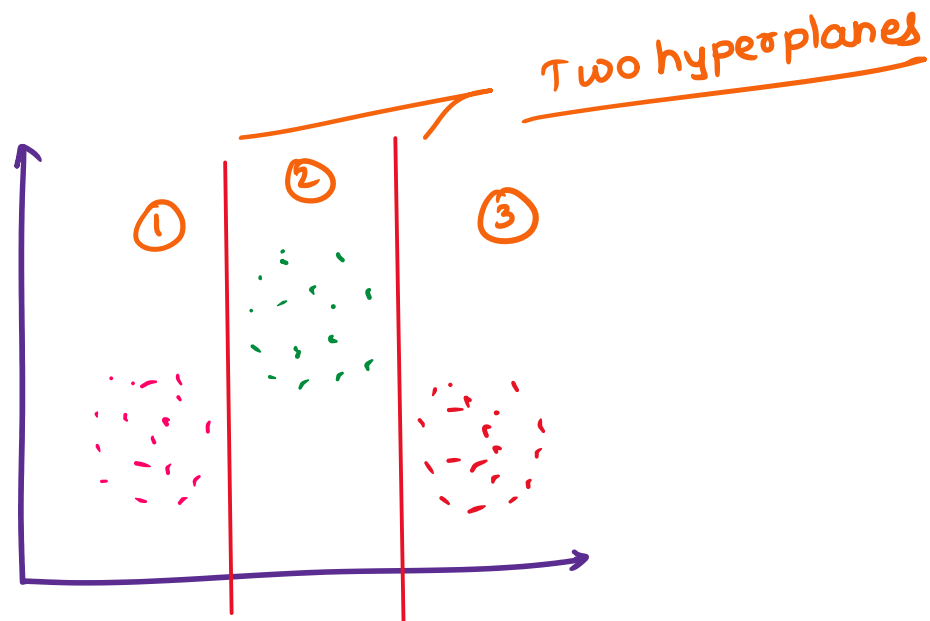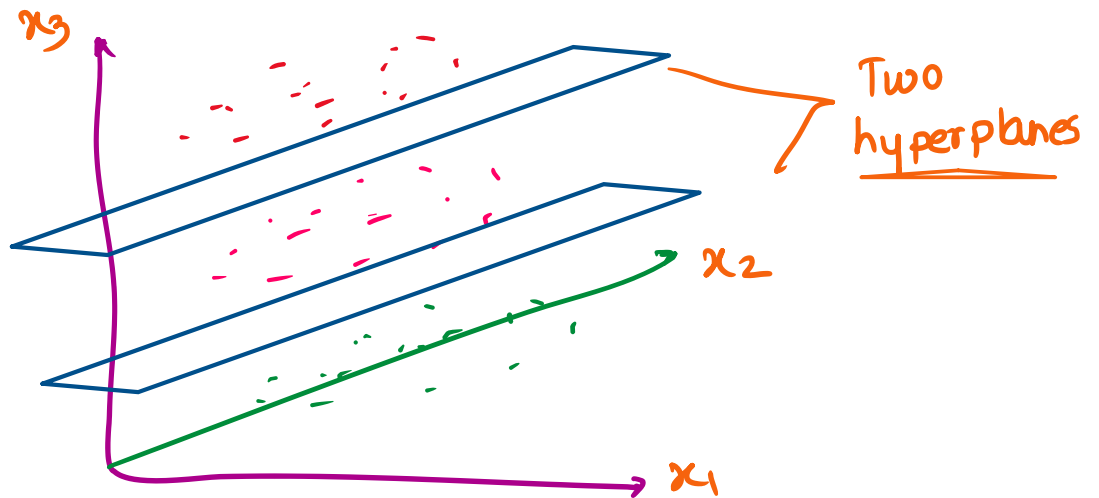
④
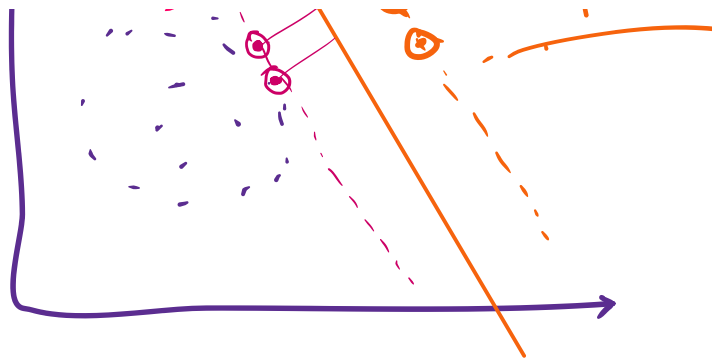
Goal of SVM → decision boundary

✓ 'Hyperplane'

To find best hyperplane

Terminology in SVM

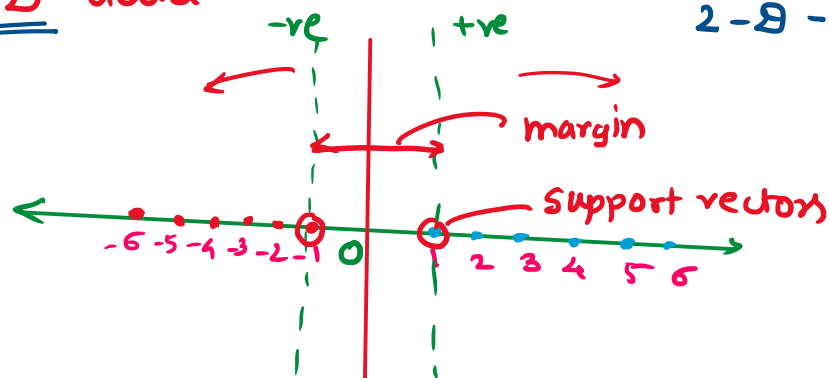2-features →
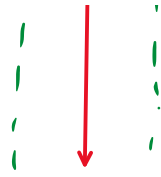1 D - hyperplane
straight line

$x_2$

-ve
hyp.

+ve
hyp.

Support
vectors

Red

support vectors

Red

support vectors

Blue

margin (max.)

Hyperplane

$x_1$

$x_3$

$x_2$

$x_1$

✓ 3-features

2D - hyperplan

4-features ——— 3-D hyperplane

① ✓② ③

Green

Pink

① - ⑮
✓② ⑰ max
③ ⑬
①

+ve

-ve

2 SV

3 SV

Two
hyperplanes

$x_3$

$x_2$

$x_1$

Two hyperplanes

① ② ③

① For **1-Ɵ** data

-ve        +ve

margin

support vectors

-6 -5 -4 -3 -2 -1   0   1   2   3   4   5   6

1-Ɵ — 1 line
2-Ɵ — x-y-plane
      ⌐y
      └x
3Ɵ - ⌐²y
       └x

② case ②

Best hyperplane



skip two missclassified data points

draw hyp. on remaining correctly classified dt.pts.

Case - 3

$y = x^2$

$2D \to$ hyp. 1D

Blue

Hyperplane

pink

X

SVM

Kernel

convert

Lower D to Higher D

$1D \to 2-D$

$-7 -6 -5 -4 -3 -2 -1 \ 0 \ 1 \ 2 \ 3 \ 4 \ 5$

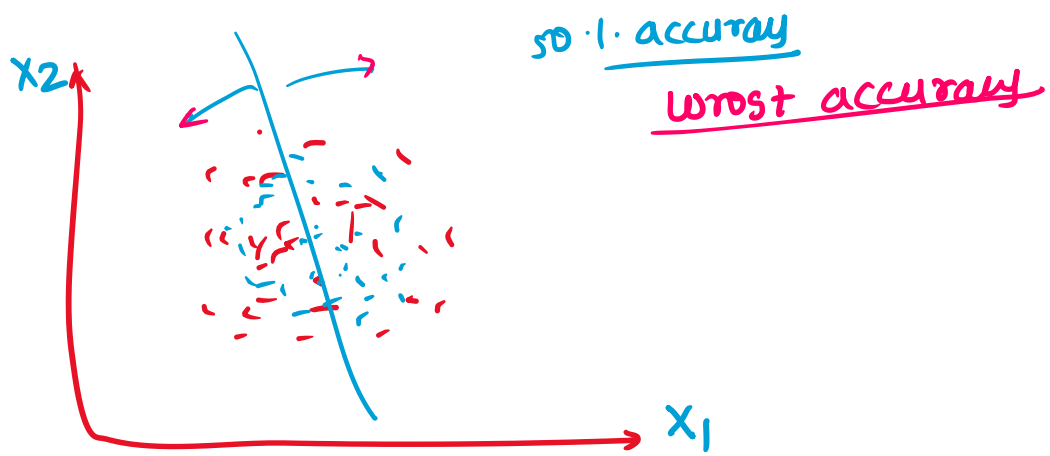Polynomial   $f(x) = x^2$   Assume $y = f(x)$

Kernel $\to$ select best function

Best kernel

by using hyperparameter tunning

Kernels  ① polynomial
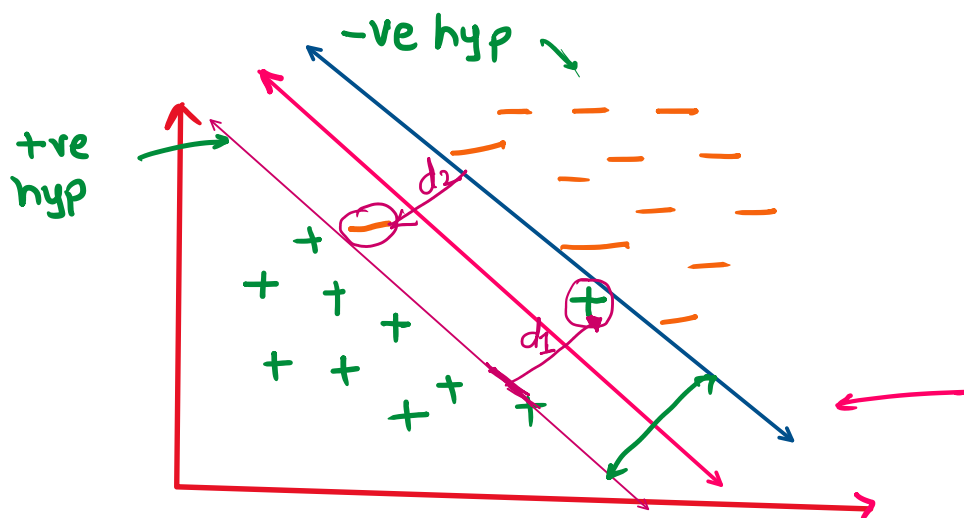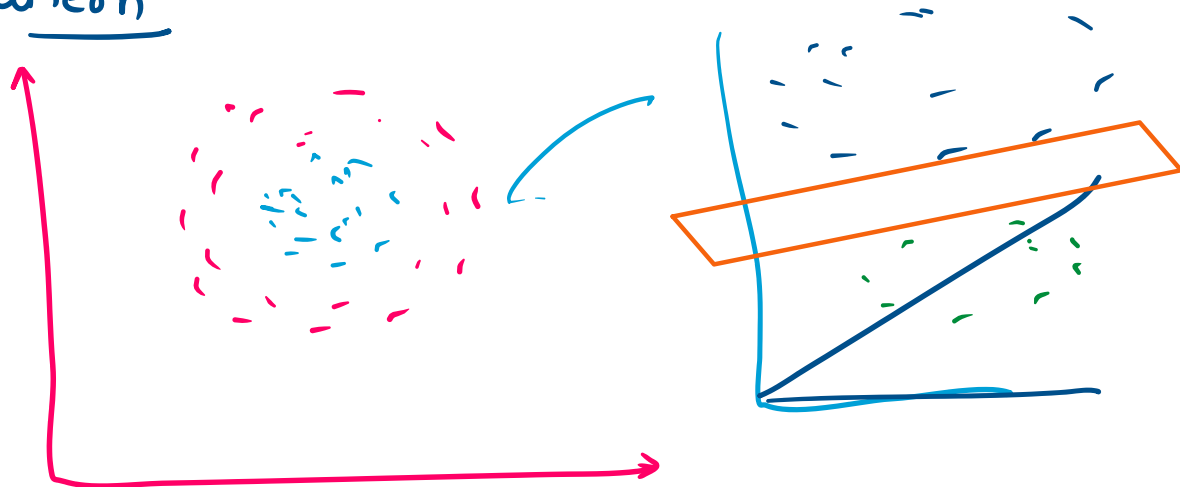
② RBF $\to$ Radial Basis Function (By default)

③ Sigmoid kernel

④ Linear kernel

X2

50·1· accuray

$X_2$

so 1% accuray

worst accuracy

$X_1$

some pattern

-ve hyp

+ve hyp

$d_2$

$d_1$

+
+  +  +
+  +  +  +
+  +  +
+

SVM Errors = Marginal Error + -Classification Error

High Margine → clear & distinct classification
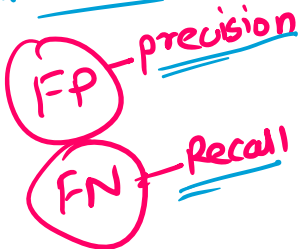less|min error
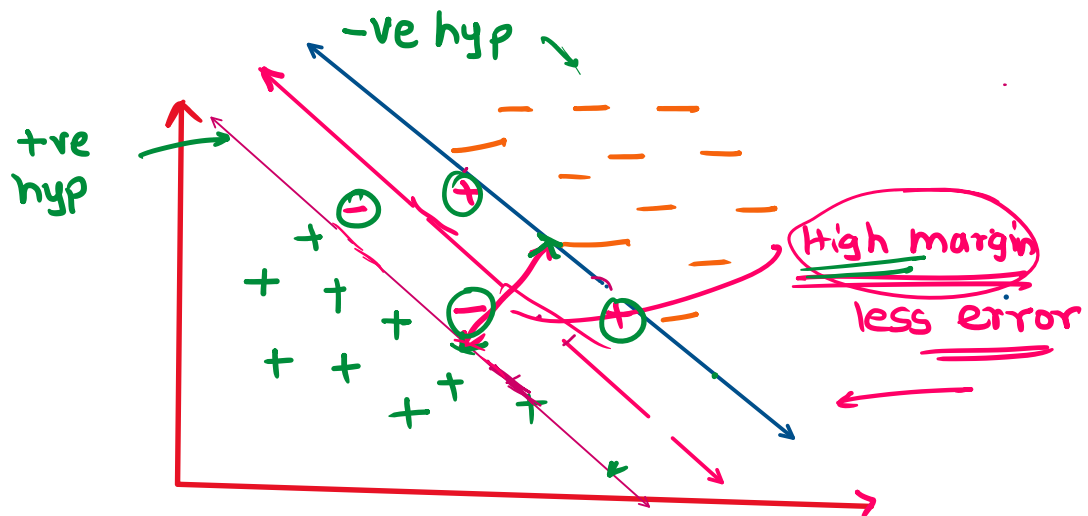
## Classification Error :-

$$\sum (d_1 + d_2)$$

$$\text{C} \cdot \sum_{i=1}^{n} \zeta \qquad \text{(zeta)}$$

SVM Errors = Marginal Error + Classification Error

Marginal Error → Concern (min)

Confusion Matrix

FP — precision
FN — Recall



+ve hyp

−ve hyp

$\ominus$ $\oplus$

High margin less error

## Soft Margin

SVM Errors = Marginal Error + $-$Classification Error

Classification Error → Concern

(min. classification Error)

less margin

& all the data points are correctly classified

# Hard Margin

$$\text{SVM Errors} = \text{Marginal Error} + \frac{-\text{Classification}}{\text{Error}}$$

(zeta)

$$\sum_{i=1}^{n} \zeta$$

c value is

$c = low$

focusing on

Too high

Focusing error   1000

Classification Error

High

Default value of   $C = 1.0$