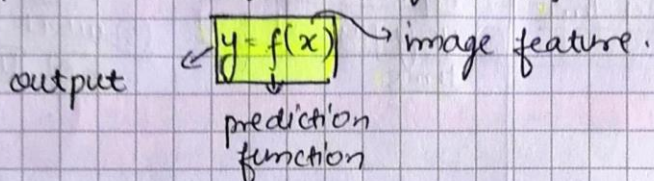## Object Recognition and SVM.

① The goal of recognition / classification is to identify in the image different description of the objects. we can have image tagging, object detection (for eg. find pedestrian) Activity recognition (usually on a sequence of images) and image parsing (reading the image) and automatically generation of text description from image. and image classification.

② The statistical learning framework suggests to apply a prediction function "$f()$" to a feature representation of the image to get the desired output.
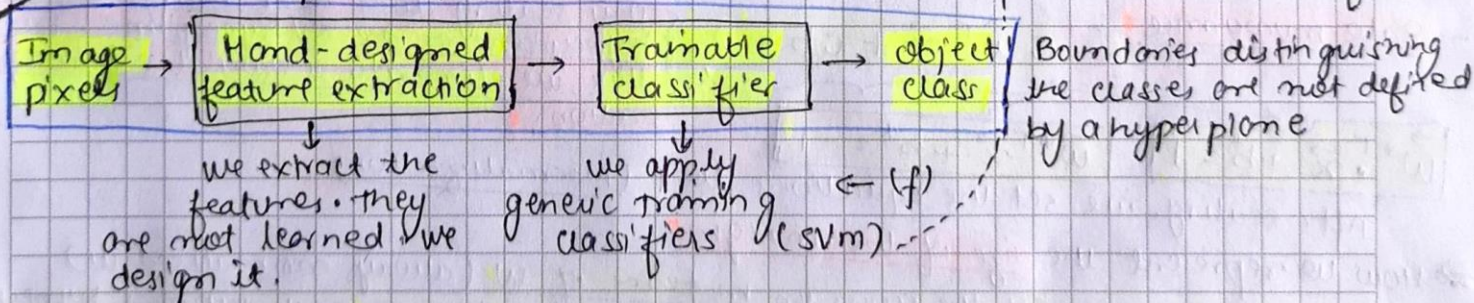
$$ y = f(x) \rightarrow \text{image feature.} $$

output ← prediction function

③ **Training** → for training, we have a training set of ^labelled examples, we know the input and output and we estimate the prediction function "$f$" by minimizing the prediction error on the training set.

**Testing** → during testing, we apply the prediction function to a new test example (eg image feature) and see the predicted value i.e. $y = f(x)$

**Steps** → ① we have a training set which consists of training images.
② we define some image features, with Training label and feed it to training model.
③ the model learns the prediction function.

**Testing** → ① we have an image never seen before. we set the image features and learned model "$f$"

④ **Traditional recognition pipeline:**      |    Non-linear classifiers

Image pixels → | Hand-designed feature extraction | → | Trainable classifier | → object class | Boundaries distinguishing the classes are not defined by a hyperplane

we extract the features. they are not learned. we design it.    we apply generic training classifiers (SVM) ← (f)

⑤ **Classifier: Nearest Neighbour** → for eg. we have 2 training examples of diff. class. and we want to see, in which category does our test example belongs. we will simply calculate the distance function for our input. we will take the smallest distance and see what label does that sample has. We will assign our point with that label.

**KNN classifier** → Instead of looking at nearest neighbour, we define a no. "$k$" which signifies how much region would be covered in that point.
eg. if $k=10$: we will take 10 closest points around our sample points.
After that we will vote for the class label and label our point with more no. of vote.

outliers can be defined as the points which lie outside of the major region. KNN is more robust to outliers. KNN is non-linear

→ **linear classifiers** → A hyperplane distinguishing the classes.

→ ① There is a hyperplane distinguishing b the planar classes.
② we see on which side does our poi lie.
③ $f(x) = \text{Sign}(w \cdot x + b)$ → offset
   ↳ our point
   +ve → one side. ⎫ different classe
   -ve → another side. ⎭

$w$ → vector     dot product = scaler
$x$ → vector  ———————→
$b$ → scaler ←

| NN advantageous | NN disadvantages | Linear class. Adv | Linear class. disa. |
|---|---|---|---|
| ① Simple to implement | ① Need good distance function. | ① Low dimensional parametric representation. | ① works for only 2 class |
| ② Decision boundaries are not necessarily linear | ② Slow at test time | ② Very fast at test time | ② we have to see here from the linear functi |
| ③ works for any no. of classes | | | ③ what if the data i not linearly separab |
| ④ Non-parametric method | | | |

→ **Support Vectors** → for every decision hyperplane, there exists a decision boundary which is calculated by maximising the margin. The input points (vectors) lying on the decision boundary are called support vectors.

→ **Decision hyperplane** can be defined as $\boxed{\vec{w}^T \cdot \vec{x} + b = 0}$

→ for taking a decision, we will use $\boxed{D(\vec{x_i}) = \text{Sign}(\vec{w}^T \cdot \vec{x} + b)}$  Decision function.
and look at the sign.

→ **Margin Hyperplanes:** $\vec{w}^T \cdot \vec{x} + b = r$
$\vec{w}^T \cdot \vec{x} + b = -r$

→ **Scale invariance** → $\boxed{C\vec{w}^T \cdot \vec{x} + Cb = 0}$ ⎫ This scaling does not change the decisio hyperplane or the SV hyperplane. But wit this, we eliminate a variable from optimi zation.

$\boxed{\vec{w}^{*T} \cdot \vec{x} + b^* = 1}$  $\boxed{\vec{w}^{*T} \vec{x} + b^* = -1}$
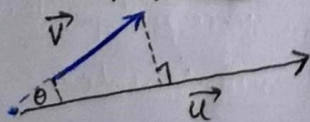
After scaling, we set the equations to 1 (Normalization).

⇒ How to represent the size of the margin.
→ we represent the size of the margin in terms of w (always normalised w) by taking that there must be atleast one point that lies on each support hyperplane so that $\boxed{w^* \cdot x_1 + b = 1}$ and $\boxed{w^T \cdot (x_1 - x_2) = 2}$
$\boxed{w^T \cdot x_2 + b = -1}$

⇒ The vector $w^*$ is perpendicular to the decision hyperplane i.e. if the dot product of 2 vectors (w and $x_1 x_2$) equals zero, the two vectors are perpendicular. The length of this vector (w) will be the distance b/w 2 hyperplanes.

• Vector projection : projection of $\vec{V}$ into $\vec{u}$   $\vec{V} \cdot \vec{u} = \|\vec{V}\| \|\vec{u}\| \cos\theta$

$\vec{V}$
$\theta$
$\vec{u}$

$\cos\theta = \dfrac{\text{adjacent}}{\text{hypotenuse}} \dfrac{\|\vec{V}\|}{\|goal\|}$   $\cos\theta = \dfrac{\|goal\|}{\|\vec{V}\|}$ . $\dfrac{\|\vec{V}\| \|\vec{u}\|}{\|\vec{V}\| \|\vec{u}\|}$ $\cos\theta = \dfrac{\|goal\|}{\|\vec{V}\|}$

$$\frac{\vec{v} \cdot \vec{u}}{\|\vec{v}\| \|\vec{u}\|} = \frac{\|goal\|}{\|\vec{v}\|} \Rightarrow \boxed{\frac{\vec{v} \cdot \vec{u}}{\|\vec{u}\|} = \|goal\|}$$

→ i.e. the margin is the projection of $x_1 - x_2$ into $w$ (the normal of the hyper plane)

$$\boxed{w^T (x_1 - x_2) = 2}$$

→ Projection → $\frac{\vec{v} \cdot \vec{u}}{\|\vec{u}\|} \vec{u} \simeq \frac{\vec{w}^T (x_1 - x_2)}{\|\vec{w}\|} \cancel{\vec{w}} \simeq$ Size of the margin $= \boxed{\frac{2}{\|\vec{w}\|}}$

→ **Goal: maximize the margin :** $\boxed{max \frac{2}{\|\vec{w}\|}}$ OR $\boxed{min \|\vec{w}\|}$

With this, we can have a decision boundary which separates the regions in a linear fashion where $t_i (\vec{w}^T x_i + b) \geq 1$

$w^T x_i + b \geq 1$ if $t_i = 1$

$w^T x_i + b \leq -1$ if $t_i = -1$

→ A problem : We want to have the best Hyperplane i.e. $\underline{min \|\vec{w}\|}$ where.
$\underbrace{t_i (\vec{w}^T x_i + b) \geq 1}_{constraint}$ .

usually, the optimisation is done with "Lagrange multipliers". **minimization**

*Imp*

$$\boxed{L(\vec{w}, b) = \frac{1}{2} \vec{w} \cdot \vec{w} - \sum_{i=0}^{N-1} a_i [t_i ((\vec{w} \cdot \vec{x_i}) + b) - 1]}$$ Main equation

Optimize Primal → Lagrange multiplier + KKT condition

* If we show or optimize the Primal, we can solve the minimization problem.

→ optimising the primal : Partial differentiation wrt b $\frac{\partial L(\vec{w}, b)}{\partial b} = 0$

Partial diff $\frac{\partial L(\vec{w}, b)}{\partial \vec{w}} = 0$
wrt w

$\sum_{i=0}^{N-1} a_i t_i = 0$

$\vec{w} - \sum_{i=0}^{N-1} a_i t_i x_i = 0 ; \quad \vec{w} = \sum_{i=0}^{N-1} a_i t_i x_i$

Substituting the value of $\vec{w}$ in the Primal equation we get :

*Imp* 🐾

$$\boxed{W(a) = \sum_{i=0}^{N-1} a_i - \frac{1}{2} \sum_{i,j=0}^{N-1} a_i a_j t_i t_j (\vec{x_i} \cdot \vec{x_j})}$$ where $a_i \geq 0$

we solve this quadratic ~~equation~~ Programming. to identify the Lagrange's multipliers

KKT condition → constraint & Lagrange multiplier = 0
$$\boxed{a_i (1 - t_i (\vec{w}^T x_i + b)) = 0}$$

if $a_i \neq 0$, then $t_i (\vec{w}^T x_i + b) = 1$ i.e. only points in the decision boundary contributes to the solution.

Final decision function
$$D(\vec{x}) = sign(\vec{w}^T \vec{x} + b) = sign\left(\left[\sum_{i=0}^{N-1} a_i t_i \vec{x_i}\right]^T \vec{x} + b\right) = sign\left(\left[\sum_{a=0}^{N-1} a_i \varepsilon_i (\vec{x_i}^T \cdot \vec{x})\right] + b\right)$$

**SVM** → It is based on the idea that the original input space can alwa[ys] be mapped to some higher dimensional feature space where the training set is separable.

**Kernel trick** → It states that, instead of explicitly computing the lifting transformati[on] $\phi(x)$, define a kernel function K such that $K(x,y) = \phi(x) \cdot \phi(y)$

→ (To be valid, the kernel function must satisfy mercer's condition)

$$[wx + b = \sum_i a_i y_i x_i \cdot x + b]\ \text{original function linear svm.}$$

↳ learned weight  ↳ support vector

**kernel trick**
$$\sum_i a_i y_i \phi(x_i) \cdot \phi(x) + b = \sum_i a_i y_i K(x_i, x) + b$$

this gives a non linear decision boundary in th[e] original feature space.

## Different kernels

① Polynomial kernels → $k(x,y) = (c + x \cdot y)^d$

② Gaussian " " → $k(x,y) = \exp\left(-\frac{1}{\sigma^2}\|x-y\|^2\right)$

## Support vector machines

**Pros:**
① Kernel based framework is very powerful, flexible
② Training is convex optimization, global optimal solution can be found
③ Amenable to theoretical analysis
④ Works very well, even with very small training samples.

**Cons:**
① There is no direct multi-class svm, must combine 2 class svm
② computation, memory.

Training Error } underfitting → Training and test error are both high (model is too simple)
Testing Error   overfitting → Training error is low but testing error is high (model is too complex)