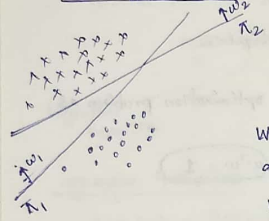## SVM (Support Vector Machines)

### Geometric Intuition
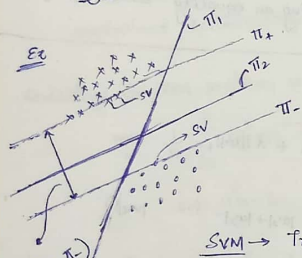


There will be many possible hyperplanes that separate positive $(+)$ve pts from $(-)$ve pts.

We had learnt in LR that the probability of a point belonging to any class given it is very close to the hyperplane will be close to $\boxed{0.5}$

So, we want a hyperplane that separates $(+)$ve pts and $(-)$ve pts as far away as possible.

→ such a hyperplane is called **margin-maximizing hyperplane.**

### Key Idea of SVM



$\Pi_2$ better than $\Pi_1$.

$\Pi_+ \rightarrow$ parallel to $\Pi$ & touches the closest $(+)$ve point to $\Pi$.

$\Pi_- \rightarrow$ parallel to $\Pi$ & touches the closest $(-)$ve point to $\Pi$.

$\underline{SVM} \rightarrow$ Try to find a $\Pi$ that maximizes the margin $(dist(\Pi_+, \Pi_-))$.

dist$(\Pi_+, \Pi_-)$ $=$ margin.

* The pts through which $\Pi_+$ & $\Pi_-$ pass through are called support vectors.

### Alternative geometric formulation of SVM

① perform convex-hull of $(+)$ve & $(-)$ve pts.
② find the shortest line connecting these hulls.
③ Bisect it.



---

### Mathematical formulation

We need to find a hyperplane.



$\Pi :-$ margin maximizing.

$\Pi : w^T x + b = 0$  ($w$ not necessarily unit vector).

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

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

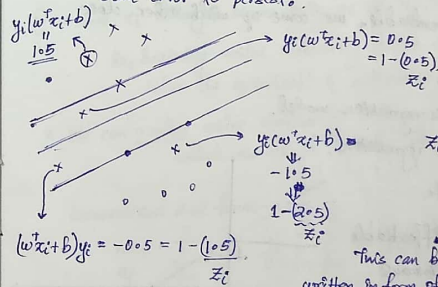By simple co-ordinate concepts we can get

$$d = \frac{2}{\|w\|_2}$$

So, we have to maximize '$d$' with the constraint that all $(+)$ve pts lie above $\Pi_+$ & all $(-)$ve pts below $\Pi_-$.

So, the optimization function becomes

$$\text{constraint optimiz}^n \text{ problem}\begin{cases}(w^*, b^*) = \underset{(w, b)}{argmax}\ \dfrac{2}{\|w\|} \quad\text{---- margin}\\[2mm] \text{s.t.}\ \ y_i(w^T x_i + b) \geq 1\ \ \forall x_i \in D_{train}\\ \phantom{s.t.\ \ }y_i\end{cases}\ \begin{array}{l}\text{Hard-margin}\\ \text{SVM constraint}\end{array}$$

But what if we have a dataset that's not linearly separable but almost linearly separable?

→ We formulate an alternative **soft-margin** constraint that allows some error to persist.



$y_i(w^T x_i + b) = 0.5$
$\quad\quad = 1 - (0.5)$
$\quad\quad \quad z_i$

$y_i(w^T x_i + b) =$
$\quad\quad - 0.5$
$\quad\quad\ \downarrow$
$1 - (1.5)$
$\quad z_i$

$(w^T x_i + b)y_i = -0.5 = 1 - (1.5)$
$\quad\quad\quad\quad\quad\quad\quad\quad z_i$

This can be written in form of eq$^n$ as
$$\boxed{\begin{array}{c}y_i(w^T x_i + b) \geq 1 - z_i\\ z_i \geq 0\end{array}}$$

We define a new metric $z_i$ for all pts such that

$$z_i = \begin{cases}0,\ \text{for correctly classified}\\ \quad\ \text{i.e } y_i(w^T x_i + b) \geq 1.\\[2mm] z_i > 0\ \&\ \text{proportional to dist away fm the correct hyperplane} (\Pi_+\text{ or }\Pi_-).\\ \quad\ \text{, for incorrectly classified pts.}\end{cases}$$

$\boxed{\begin{array}{l}\text{correctly \&}\\ \text{incorrectly wrt } \Pi_+ \,\&\, \Pi_-\\ \text{considered here.}\end{array}}$

Now, we want to minimize these errors.

So, our optimization function becomes

$$\underset{w, b}{\arg\min} \left(\frac{\|w\|}{2}\right) + c\left(\frac{1}{n}\sum_{i=1}^{n} z_i\right)$$

maximizing $\frac{2}{\|w\|}$ same as

minimizing $\frac{\|w\|}{2}$

$\hookrightarrow$ avg cost of misclassified pts.

(regularzn) $\leftarrow \frac{1}{\text{margin}}$

s.t. $y_i(w^T x_i + b) \geq 1 - z_i \quad \forall i$ ← hinge-loss

def$^n$ of $z_i$ → $z_i \geq 0$

This is soft-margin formulation of SVM.

$c \rightarrow$ hyper-parameter.

$c\uparrow \Rightarrow$ tendency to make mistakes $\downarrow \Rightarrow$ overfit $\Rightarrow$ high variance. on $D_{train}$

$c\downarrow \Rightarrow$ underfit $\Rightarrow$ high bias.

$$\boxed{c = \frac{1}{\lambda}}$$
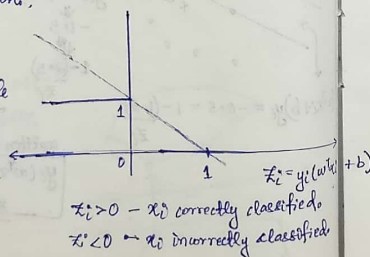
## Loss-minimization interpretation

We earlier saw that $0-1$ loss is ideal ~~loss function~~. But due to it being non-differentiable, we come up with other close approximate loss functions.

logistic loss gave $\longrightarrow$ Logistic regression, model.

squared loss gave $\longrightarrow$ linear regression.

hinge loss gives $\longrightarrow$ SVM.

* Although hinge loss also isn't diffentiable at $z_i = 1$, but since it is continuous unlike $0-1$ loss, we can work around hacks to work with it.

$z_i = y_i(w^T x_i + b)$

$z_i > 0$ → $x_i$ correctly classified.
$z_i < 0$ → $x_i$ incorrectly classified

hinge-loss :– $0, \quad z_i \geq 1$. Alternatively,
$\qquad\qquad 1 - z_i, \quad z_i < 1$ $\quad$ hinge-loss $= \max(0, 1 - z_i)$

## Dual form of SVM

$$\begin{cases} \underset{w,b}{\min} \quad \frac{1}{2}\|w\| + c\sum_{i=1}^{n} z_i \\ \\ \text{s.t. } y_i(w^T x_i + b) \geq 1 - z_i \; \forall i \\ \qquad \& \; z_i \geq 0. \end{cases}$$

(Primal form) of Soft-margin SVM.

$$\underset{\alpha_i}{\max} \sum_{i=1}^{n} \alpha_i - \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \; x_i^T x_j$$

s.t. $\alpha_i \geq 0$

$\& \; \sum_{i=1}^{n} \alpha_i y_i = 0$

(Dual form) of Soft margin SVM.

Important observations from dual-form formulation.

① for every $x_i$, we have an $\alpha_i$

② $x_i$s only occur in the form of $x_i^T x_j$.

On solving we get the model as

$$\boxed{f(x_q) = \sum_{l=1}^{n} \alpha_i y_i \, x_i^T x_q + b}$$

→ for non SV, $\alpha_i = 0$. So, for $f(x_q) \rightarrow$ only the support vectors matter.

③ $\alpha_i > 0$ only for support vectors, else $0$.

④ – since $x_i$ always occurs only as $\boxed{x_i^T x_j} \Rightarrow x_i \cdot x_j \Rightarrow$ cosine_sim $(x_i, x_j)$.

So, basically cosine similarity values are only required to solve the optimizat$^n$ & ultimately get the model.

* We can replace cosine-similarity with any other kind of similarity, which makes it much more powerful.

Generalized dual form →

$$\underset{\alpha_i}{\max} \sum_{i=1}^{n} \alpha_i - \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n} \alpha_i \alpha_j \, y_i \, y_j \; K(x_i, x_j)$$

← called kernel function can be any kind of similarity b/w $x_i$ & $x_j$.

Even during evaluation $x_i$ occurs only as $x_i^T x_q$

$$f(x) = \sum_{i=1}^{n} \alpha_i y_i \, x_i^T x_q + b$$

$\hookrightarrow K(x_i, x_q)$

So, we can use similarity values for even prediction as well.

## Kernel trick

Intuitively if we see, SVM & LR are very similar.

* The most important idea of SVM is the kernel trick.

Kernel trick → replacing $x_i^T x_j$ with a generalized similarity function $K(x_i, x_j)$

This gives SVM potential to perform in non-linear datasets.



linear-SVM — fail.
log-regn" — fail.
log-regn" + feature transform — succeed.
kernel SVM + right kernel — succeed.

Just like in feature transformatn, we map data pts to a new space, similarly, the kernel in SVM does the mapping & then tries to find a ~~plan~~ hyperplane in that mapped ~~by~~ space.

## Polynomial Kernel

Given circular data like above, we can

① handle using feature transforms.

$$(f_1, f_2) \xrightarrow{F \cdot T \cdot} (f_1^2, f_2^2) — \boxtimes$$

② fit right kernel.

General polynomial kernel of form —

$$\boxed{K(x_1, x_2) = (x_1^T x_2 + c)^d}$$

A quadratic kernel can be written as.

$$K(x_1, x_2) = (x_1^T x_2 + 1)^2 \qquad \text{assuming } c = 1.$$

let $x_1 \; \&< x_{11}, x_{12} >$ & $x_2 < x_{21}, x_{22} >$

---

$$K(x_1, x_2) = (1 + x_{11} x_{21} + x_{12} x_{22})^2$$

$$= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 + 2 x_{11} x_{21} + 2 x_{12} x_{22} + 2 x_{11} x_{21} x_{12} x_{22}$$

Say $< 1, x_{11}^2, x_{12}^2, \sqrt{2} x_{11}, \sqrt{2} x_{12}, \sqrt{2} x_{11} x_{12} > — x_1'$

$< 1, x_{21}^2, x_{22}^2, \sqrt{2} x_{21}, \sqrt{2} x_{22}, \sqrt{2} x_{21} x_{22} > — x_2'$

Then the product we have got due to quadratic kernel is equivalent to

$$x_1'^T x_2' \, .$$

* So, the kernel basically internally 1st maps $x_1$ & $x_2$ to $x_1'$ & $x_2'$ and then finds their similarity (cosine).

$$\text{kernelization} \to d \xrightarrow[\text{internally}]{\text{feature transform}} d' \qquad (d' > d)$$

* In our example we went fm 2-D data to 6-D data.

That is what Mercer's theorem says.

It says $\underline{Kernel\ trick}$ is basically nothing but

$$d \xrightarrow[\text{transform}]{\text{implicit}} d' \qquad (d' > d).$$
↳ data linearly separable in this space.

The challenge is to select the right kernel.

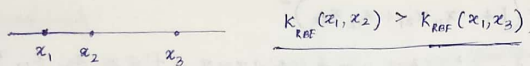## RBF - Kernel (Radial Basis Function)
↓
most popular and general purpose kernel.

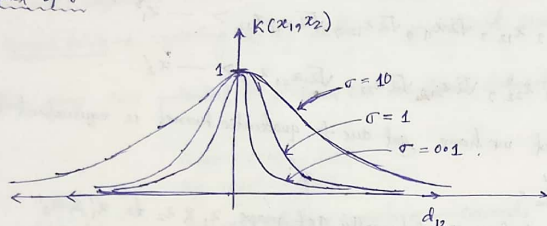$$\boxed{K_{rbf}(x_1, x_2) = \exp\left\{ \frac{-\|x_1 - x_2\|^2}{2\sigma^2} \right\}}$$

hyperparameter.

$$\boxed{K(x_1, x_2) = \frac{1}{e^{d_{12}^2/2\sigma^2}}}$$

$d_{12} \uparrow, \quad K(x_1, x_2) \downarrow.$

$x_1 \quad x_2 \qquad x_3$

$$K_{RBF}(x_1, x_2) > K_{RBF}(x_1, x_3)$$

**Impact of $\sigma$**

$K(x_1, x_2)$

1

$\sigma = 10$

$\sigma = 1$

$\sigma = 0.1$

$d_{12}$

for ~~large~~ small $\sigma$, the kernel value (similarity value) decreases very quickly with increasing distance.

Whereas for large $\sigma$, the kernel value decreases gradually, slowly.

So, small $\sigma$ is similar to small 'k' in KNN. — considering very few closest neighbours.

large $\sigma \implies$ large k $\rightarrow$ considering more neighbours which are even a bit farther.

K-NN works pretty good but its biggest drawback is runtime ~~so~~ complexity
$\hookrightarrow$ needs to store all the datapoints.
or if optimised, LSH.

In RBF-SVM, we just need to store $\alpha_i$'s of SVs.
That's it.
And usually, #SVs $<< n$

If can't decide which kernel, go with RBF.
2 hyperparameters — $c, \sigma$. —— Random search.

---

**Domain Specific Kernels**

* With research over the years we have come up with domain specific kernels.

  Ex $\rightarrow$ String kernels for text classification,
  genome kernels (Bioinformatics),
  graph kernels (for graph-data).

  Domain knowledge + Appropriate feature transformations $\longrightarrow$ Custom-domain specific kernels.

---

**Train and Runtime Complexity**

Train:— SGD (stochastic Gradient Descent),
Specialised algo to solve dual form $\rightarrow$ sequential minimal Optimization (SMO).

libSVM $\rightarrow$ best library for training SVM.
(c/c++)

Training time:— $O(n^2)$ for kernel SVMs.
(2007):— $O(nd^2)$ if $d < n$

If $n$ is large, $O(n^2)\uparrow\uparrow$
So, typically don't use SVM when n is large.

Runtime:→ $f(x_q) = \sum_{i=1}^{n} \alpha_i y_i k(x_i, x_q) + b$.

$\uparrow$ depends on

$\alpha_i = 0$ for non-SV

#SVs = k (say)
then runtime complexity = $O(kd)$ $\qquad 1 \le k \le n$

logistic regression- runtime — $O(d)$

In normal SVM we have no control over no. of support vectors.

## nu-SVM

c-SVM — original formulation.    c — hyperparameter  $c \geq 0$

nu-SVM — alternative formulation.  nu — hyperparameter  $0 \leq nu \leq 1$

$$\underline{nu \geq \text{fraction of errors.}}$$
$$\underline{nu \leq \text{fraction of SVs}}$$

Eg— If we want to ~~have~~ permit an error of 10% in worst case,
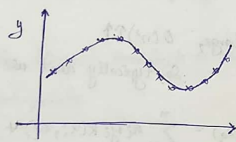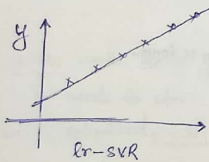
$$nu = 0.1$$

· Also, #SVs is greater than 10% of $\underline{n}$.

We can have an upper bound on error, but not on $\underline{\text{#SVs.}}$

---

## Support Vector Regression (SVR)

Maths:—  $\min\limits_{w,b} \dfrac{1}{2}\|w\|^2$    s.to  $\{y_i - (w^T z_i + b)\} \leq \epsilon$

&  $(w^T z_i + b) - y_i \leq \epsilon$      $\epsilon \geq 0$
                                              $\downarrow$
                                     hyper-parameter.

· This can also be kernelised.



lr—SVR                    kernel— SVR.

$\epsilon \downarrow \Rightarrow$ errors low on training data $\Rightarrow$ overfit $\uparrow \Rightarrow$ high variance

$\epsilon \uparrow \Rightarrow$ errors on $D_{train}$ higher $\Rightarrow$ underfit $\Rightarrow$ high bias.

---

## Real world cases

feature engg & transformation $\longrightarrow$ kernel design, finding the right kernel.

decision surface $\rightarrow$ lr. SVM $\rightarrow$ hyperplane.

  kernel SVM $\rightarrow$  $d$  $x_i$  $\longrightarrow$ non-linear surface
                                    $\downarrow$
                                  $d'$ $x_i'$  $\longrightarrow$ linear surface.

similarity matrix ↙

Interpretability and feature importance —
                  for lr-SVM, same as LoR.
                  for kernel SVM, difficult.

Outliers $\rightarrow$ generally v. little impact.

exception ×× RBF with small $\sigma$ = KNN with small k. —

Bias-variance $\rightarrow$  $c \uparrow \Rightarrow$ overfit
                            $c \downarrow \Rightarrow$ underfit.

large $d$ $\rightarrow$ v.good for SVMs.   as SVMs eventually map to higher
                                          dimensions only.

Best case $\rightarrow$ have right kernel ↙

worst case $\rightarrow$ $n$ is large. $\longrightarrow$ high train time
                      #SVs large. $\longrightarrow$ low latency not possible.