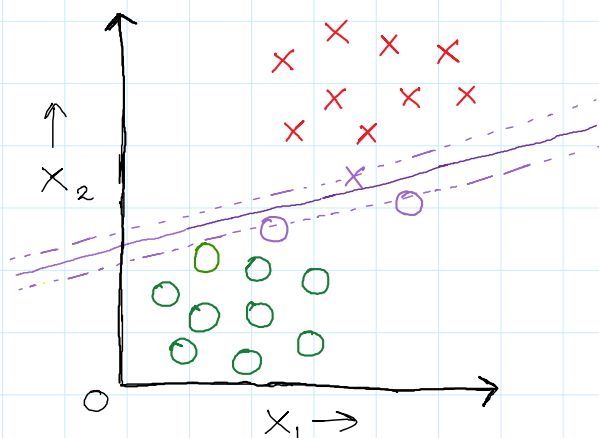


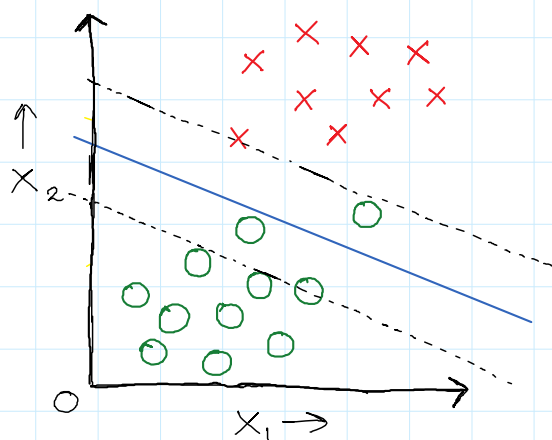
SOFT MARGIN CLASSIFICATION

For soft margin classifier, we need to allow some misclassifications so that we have a tradeoff between the margin width and the 'correctness' of our classifier. This indicates that we need to have some 'slack' or 'dheel' (दोल) in classifying all the data points correctly so that our margin width does not suffer because of noise in our dataset.

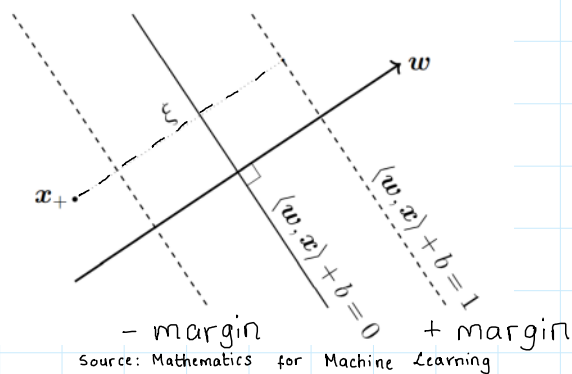
Eg: It is unlikely that any real world data will be linearly separable, even if it is, it might look somewhat like :



A soft margin SVM will/can allow some slackness in the process which will give us a wider margin at the cost of a few misclassifications.



So we aim at minimizing the total amounts of slacks while maximizing the margin width. Let's introduce a slack variable ξ_i for every training example $i = 1, \dots, m$.



Where ϵ quantifies the 'error' or deviation of a + sample x_+ from the +margin: $w \cdot x + b = +1$. For a SV, $\epsilon_{sv} > 1$ means that it has been misclassified as it lies on the other side of the hyperplane. If $\epsilon_{sv} < 1$, then it lies inside that margin yet correctly classified.

OPTIMIZATION:

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \epsilon_i$$

subject to ① $y_i(w \cdot x_i + b) \geq (1 - \epsilon_i) \quad \forall i = 1, 2, \dots, m$

The amount of slack for the i^{th} sample.

② $\epsilon_i \geq 0 \quad \forall i = 1, 2, \dots, m$

If C is large, then even a very little slack will increase the cost. If C is kept small, then relatively more slack can be allowed.

Thus, reducing C helps in increasing the no. of SVs inside and on the margin which results in decreasing the variance and reducing over-fitting because our model will generalize better.

New Primal Problem :

New Primal Problem :

$$\mathcal{L}(\omega, b, \xi, \alpha, \mu) = \frac{1}{2} \omega \cdot \omega + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(\omega \cdot x_i + b) - (1 - \xi_i)] - \sum_{i=1}^m \mu_i \xi_i$$

$$\alpha_i \geq 0 \\ \mu_i \geq 0 \quad \forall i = 1, 2, \dots, m$$

$$\min_{\omega, b, \xi} \quad \max_{\alpha, \mu} \quad \mathcal{L}(\omega, b, \xi, \alpha, \mu)$$

NOTE: $y_i(\omega \cdot x_i + b) \geq 1 - \xi_i$

if $\xi_i = 0$, then this i^{th} sample is correctly classified.

$$\text{if } y_i(\omega \cdot x_i + b) = 1 - 0$$

then it's a correctly classified S.V.

else if $\xi_i > 0$, then this sample might be misclassified.

$$\text{if } y_i(\omega \cdot x_i + b) = 1 - \xi_i \quad \text{and } \xi_i > 1$$

then this 'is' a misclassified data point.

else if $y_i(\omega \cdot x_i + b) \geq 1 - \xi_i$ is smaller than 1, then it's misclassified.

Partially differentiating the Primal Lagrangian wrt the Primal Variables ω, b and ξ gives us:

$$\frac{\partial \mathcal{L}}{\partial \omega} = \omega - \sum_{i=1}^m \alpha_i y_i x_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{i=1}^m \alpha_i y_i$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = C - \alpha_i - \mu_i$$

Equating them to zero gives :

$$\boxed{\omega = \sum_{i=1}^m \alpha_i y_i x_i} \quad - \textcircled{1}$$

$$w = \sum_{i=1}^m \alpha_i y_i x_i$$

- ①

$$\sum_{i=1}^m \alpha_i y_i = 0$$

- ②

$$C - \alpha_i - u_i = 0$$

- ③

Substituting ①, ② and ③ after simplifying our Primal Problem :

$$\begin{aligned} \mathcal{L}_P : & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^m \alpha_i y_i \left(\sum_{j=1}^m \alpha_j y_j x_j \right) \cdot x_i - b \sum_{i=1}^m \alpha_i y_i \\ & + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m u_i \xi_i + C \sum_{i=1}^m \xi_i \end{aligned}$$

$$\mathcal{L}_D : -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^m \alpha_i + \sum_{i=1}^m (C - \alpha_i - u_i) \xi_i$$

$$\mathcal{L}_D(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

Clearly, $\alpha_i \geq 0$ and $u_i \geq 0 \quad \forall i = 1, 2, \dots, m$ as they are Lagrangian Multipliers.

Given that $\alpha_i = C - u_i$ (from ③)

$$\Rightarrow 0 \leq \alpha_i \leq C$$

All data points for which $0 < \alpha_i < C$ holds true, are our Support Vectors that lie exactly on our margin. This is explained by Karush - Kuhn Tucker conditions.

Now, solving the Dual Problem shown above, we will get our Lagrangian Multipliers back. This enables us to get back our Primal Variables w and b which parametrize our separating hyperplane.

$$; \quad \frac{1}{m} \quad ; \quad , \quad \backslash$$

parametrize our separating hyperplane.

$$\omega^* = \sum_{i=1}^m \alpha_i y_i x_i \quad (\text{From } \textcircled{1})$$

Once we get ω , we can plug in that into:

$$\omega \cdot x_{sv} + b^* = y_{sv} \quad - \textcircled{4}$$

for some SV lying exactly on our margin.

$$\Rightarrow b^* = y_{sv} - \omega \cdot x_{sv}$$

If everything was hunky-dory, we were done at this point.

BUT

It is possible that none of the SVs lie exactly on the margin i.e. all SVs violate the margin.

Then we can't use $\textcircled{4}$ as our assumption fails.

By some more math, we can say that in such a case we can find b^* as follows:

$$b^* = \text{Median}(|y_s - \omega^* \cdot x_s|) \quad \forall s \in \text{SVs.}$$