

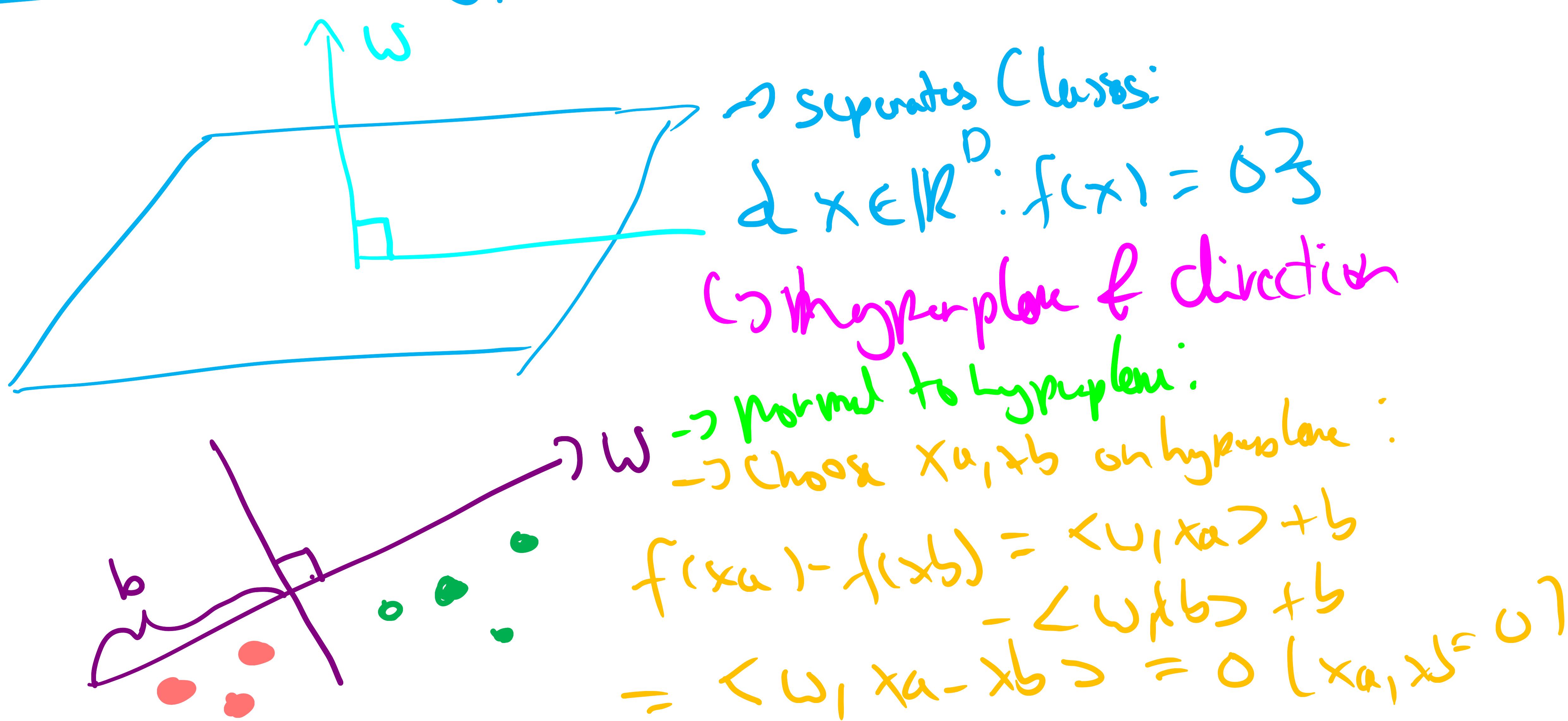
12 Support Vector Machines

- used for classification / regression
- $y_n \in \{-1, 1\}$
- linear model is considered, transforming ϕ for nonlinearity
- geometric way to think about SVM
- relies on inner products, projections
- different to Max. Likelihood:
 - Max. Likelihood: model based on probabilistic view of data
 - SVM: designs function to be optimized during training \rightarrow loss function
- Hyperplane (affine subspace) \Rightarrow linear separator

Ideas:

- 1) Hard margin
- 2) Soft margin
- 3) Dual SVM: convex hull of each class

12.1 Separating Hyperplanes



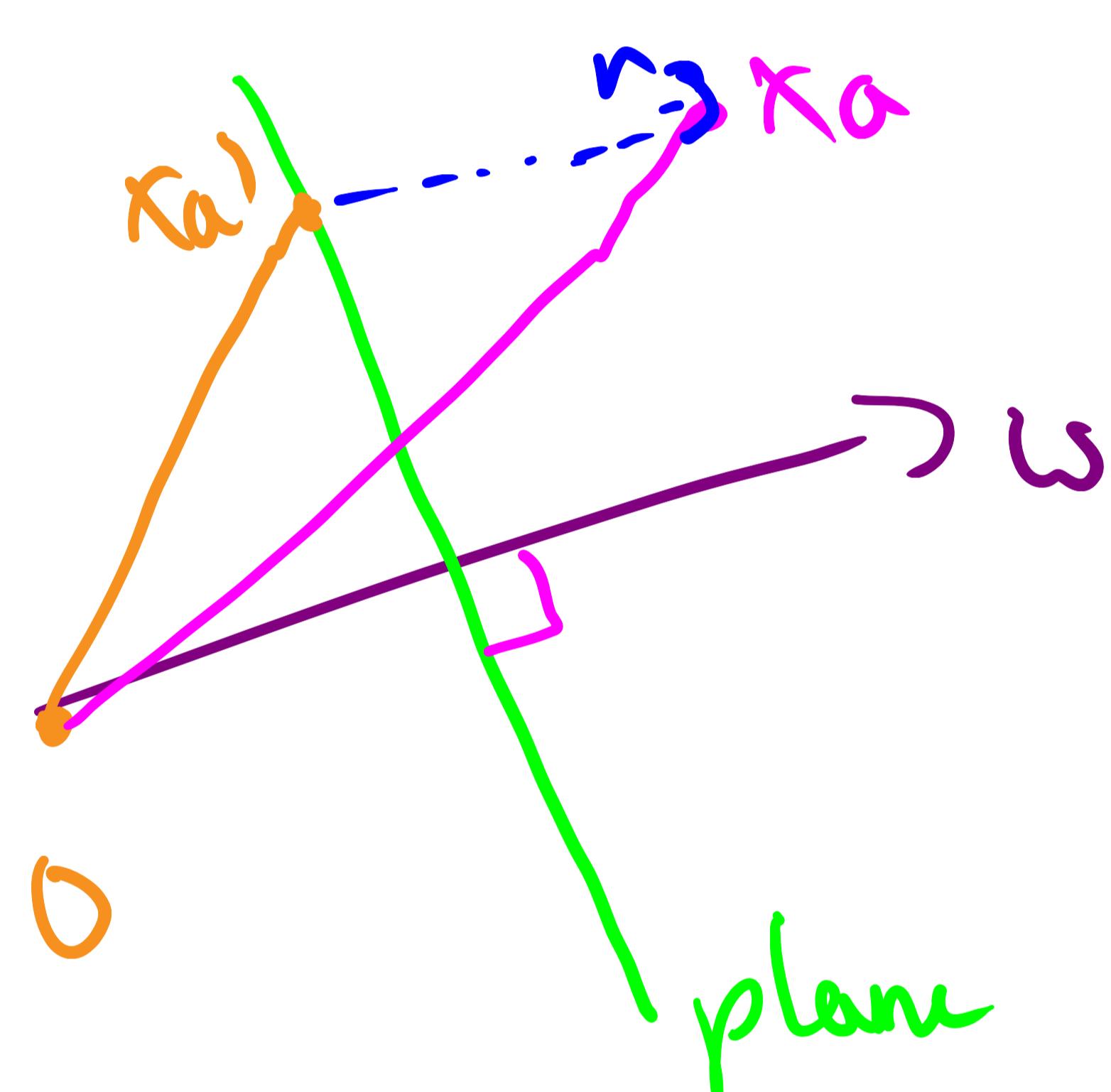
During Training

$\langle \omega, x_h \rangle + b \geq 0$, when $y_h = 1$

$\langle \omega, x_h \rangle + b < 0$, $y_h = -1$

12.2 Primal SVM

- Many hyperplanes (infinite) possible
⇒ choose plane ω max. margin
 - Key: Scale @ which d is measured?
→ scale as bound on hyperplane
- r: dist. one from hyperplane



is scaling of w :

$$x_a = x_a' + r$$

$$\frac{w}{\|w\|}$$

→ normalize by $\|w\| = \sqrt{w^T w}$

positive example:

more than r from plane

Negative example:

more than r from plane

$$y_h (\langle \omega, x_h \rangle + b) \geq r$$

SVM Constrained Optimization

Problem:

$$\max_{w/b} r$$

\uparrow
margin

subject to $y_h (\langle \omega, x_h \rangle + b) \geq r, \|w\| = 1, r > 0$

→ max r , else sum data is on right side
of the plane

\uparrow normalization

→ Var. will: large margin \Rightarrow "complexity" low

Margin Max.

$$\langle \omega, x_a \rangle + b = 0 \quad / \text{subs.}$$

$$\langle \omega, x_a - r \frac{\omega}{\|\omega\|} \rangle + b = 0 \quad / \text{bilinearity}$$

$$\langle \omega, x_a \rangle + b - r \frac{\langle \omega, \omega \rangle}{\|\omega\|} = 0$$

$$r = \frac{1}{\|\omega\|} \quad \downarrow \quad \langle \omega, \omega \rangle = \|\omega\|^2$$

$\rightarrow r$ derived last ω we do not know yet

$$\max_{\omega, b} \frac{1}{\|\omega\|}$$

$$\text{Subject to } y_n (\langle \omega_n, x_n \rangle + b) \geq 1 \quad \text{for all } n \in \{1, N\}$$

Often: squared norm is used:

$$\max_{\omega, b} \frac{1}{2} \|\omega\|^2$$

Subject to $-||-$

\rightarrow hard margin SUM (can be relaxed)

Set margin to 1

$$\max_{\omega, b} \underbrace{\underbrace{1}_{\text{margin}}}_{\text{margin}}$$

$$\text{Subject to } y_n (\langle \omega, x_n \rangle + b) \geq r, \quad \|\omega\| = 1, \quad r > 0$$

equivalent to:

$$\min_{\omega, b} \frac{1}{2} \|b\|^2$$

$$\text{Subject to } y_n (\langle \omega, x_n \rangle + b) \geq 1$$

Proof \Rightarrow book

Soft Margin SVM

→ allow for outliers

→ key idea: slack variable $\xi_n \rightarrow$ each label pair (x_n, y_n)

↳ subtract from margin

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n$$

$$\text{Subject to } y_n (\langle w, x_n \rangle + b) \geq 1 - \xi_n$$

$$\xi_n \geq 0$$

C: trades off size of margin & slack

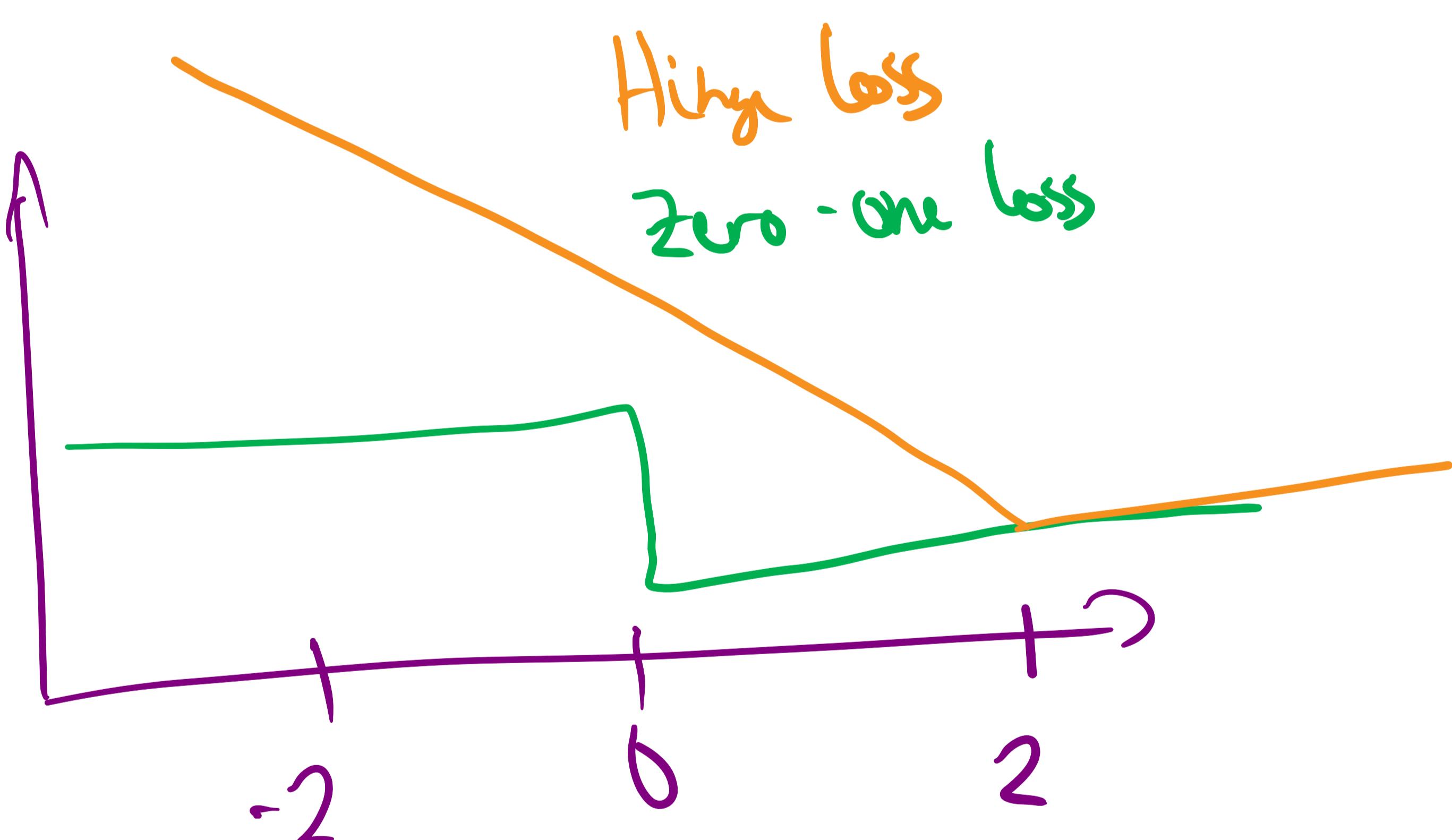
$\|w\|^2$: regularizer

$C \uparrow \rightarrow$ regularization \downarrow (slack has more weight)

Loss:

Hinge loss:

$$l(t) = \max \{0, 1-t\} \quad t = y f(x) = y(\langle w, x \rangle + b)$$



Loss:

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \max \{0, 1-y_n(\langle w, x_n \rangle + b)\}$$

$\underbrace{\|w\|^2}_{\text{regularizer}}$

"Margin max. can be interpreted as regularization"

Dual SVM

- independent on # of features (unlike primal SVM)

↳ # of params ↑ when # training examples ↑

- convex duality

- allows for kernels to be applied

$$L(\omega, b, \xi, \alpha, \gamma) = \frac{1}{2} \|\omega\|^2 + C \left(\sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n (\gamma_n (\langle \omega, x_n \rangle + b) - 1 + \xi_n) - \sum_{n=1}^N \gamma_n \xi_n \right)$$

→ partial derivatives to 0:

$$\omega = \sum_{n=1}^N \alpha_n \gamma_n x_n \quad (\text{representer theorem})$$

→ solution lies in span of the data

$$\min \alpha \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle - \sum_{i=1}^N \alpha_i$$

$$\text{Subject to } \sum_{i=1}^N y_i \alpha_i = 0 \quad 0 \leq \alpha_i \leq C \text{ for all } i=1, \dots, N$$

1) obtain α

2) get ω through representer theorem

$$\omega^* = \sum_{n=1}^N \alpha_n y_n x_n \quad (x_n \text{ on the boundary})$$

→ if no boundary: union of SVs

// Convex Hull View

• build convex set w/ examples \Rightarrow smallest set possible (convex hull)

• triangle formed by edges of region

$$\text{Conv}(X) = \left\{ \sum_{n=1}^N \alpha_n x_n \mid \sum_{n=1}^N \alpha_n = 1, \alpha_n \geq 0 \right\}$$

- do not overlap when classes are separated

1) pick points c & d in each class

$$\omega^* = c - d \rightarrow \text{require closedness}$$

$$= \arg \min_{\omega} \|\omega\|^2 = \arg \min_{\omega} \frac{1}{2} \|\omega\|^2$$

$$c = \sum_{n:y_n=+1} \alpha_n^+ x_n \rightarrow \text{must be in positive convex hull}$$

$$d = \sum_{n:y_n=-1} \alpha_n^- x_n$$

Substitute

$$\Rightarrow \min_{\alpha} \frac{1}{2} \left\| \sum_{n:y_n=+1} \alpha_n^+ x_n - \sum_{n:y_n=-1} \alpha_n^- x_n \right\|^2$$

$$\text{require: } \sum_{n:y_n=+1} \alpha_n^+ = 1 \Rightarrow \text{implies: } \sum_{n=1}^N y_n \alpha_n = 0$$

$$\sum_{n:y_n=-1} \alpha_n^- = 1$$

\Rightarrow Constrained Optimization

\Rightarrow dual hard margin SVM

To obtain soft margin: reduced hull (upper bound of convex hull)

2.4 Kernels

- dual SVM: no inner products between examples & params
- consider features $\phi(x_i)$ to represent $x_i \rightarrow$ only change inner product in dual SVM
- ⇒ classification method & representation method are separate
 $\phi(x)$ can be non-linear \Rightarrow Kernel

Kernel: function K

$K: X \times X \rightarrow \mathbb{R}$ for which there exists
a Hilbert Space H & $\phi: X \rightarrow H$ a feature map:

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle_H$$

examples

similarity measure

\Rightarrow any Kernel has unique Hilbert Space

→ Generalization from inner product \rightarrow feature map:

KERNEL TRICK