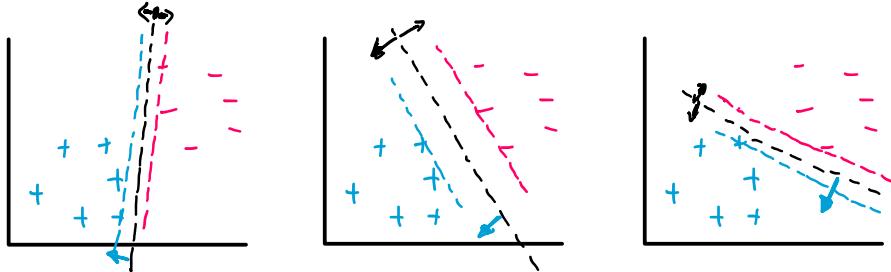


12-Nonlinear Models (SVMs)

Max-margin criteria

- Consider a binary, linearly separable classification task
Non-binary, non-linear tasks => future development
- All the below classifiers give the same training accuracy (100%), but still one looks better than others.



- It looks the classifier with the maximum margin is the most stable one. (Max-margin classification has profound theoretical justifications, but we treat it as heuristics.)

Formulation

Input features: $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$

Target: $t \in \{+1, -1\}$

Note: There's no real difference whether we call the target $\{0, 1\}$ or $\{+1, -1\}$

We would still like a linear (more rigorously, affine) model on \mathbf{x} , with parameters \mathbf{w}, b . Thus, the decision boundary is the hyperplane $\mathbf{w}^\top \mathbf{x} + b = 0$

For a data point $(\mathbf{x}^{(m)}, t^{(m)})$, the "distance" between $\mathbf{x}^{(m)}$ to the hyperplane is

$$\frac{|\mathbf{w}^\top \mathbf{x}^{(m)} + b|}{\|\mathbf{w}\|}$$

If $\mathbf{x}^{(m)}$ is in the half space that the normal vector points to, then $\mathbf{w}^\top \mathbf{x}^{(m)} + b > 0$, and vice versa.

Alternatively, we can express the distance as

$$\frac{t^{(m)}(\mathbf{w}^T \mathbf{x}^{(m)} + b)}{\|\mathbf{w}\|}$$

Thus, the training objective is to

$$\max_{\mathbf{w}} \min_m \frac{t^{(m)}(\mathbf{w}^T \mathbf{x}^{(m)} + b)}{\|\mathbf{w}\|}$$

It is noted that we may scale \mathbf{w} and b and the result is not changed.

The trick is to impose a constraint

$$\min_m t^{(m)}(\mathbf{w}^T \mathbf{x}^{(m)} + b) = 1$$

Then, the optimization is

$$\underset{\mathbf{w}, b}{\text{maximize}} \frac{1}{\|\mathbf{w}\|}$$

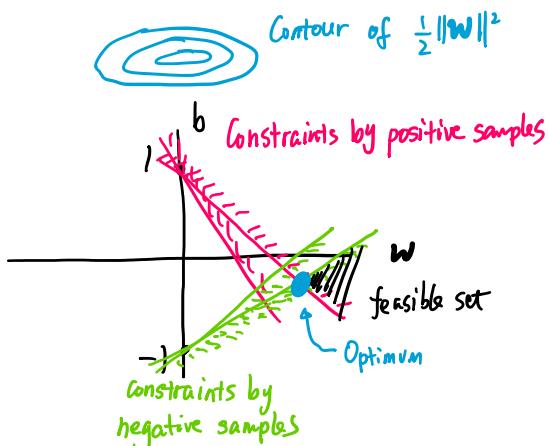
$$\text{subject to } \min_m t^{(m)}(\mathbf{w}^T \mathbf{x}^{(m)} + b) = 1$$

and is equivalent to

$$\underset{\mathbf{w}, b}{\text{minimize}} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{subject to } t^{(m)}(\mathbf{w}^T \mathbf{x}^{(m)} + b) \geq 1 \quad \text{for } m=1, \dots, M$$

as the equality constraint $\min_t t^{(m)}(\mathbf{w}^T \mathbf{x}^{(m)} + b) = 1$
will be automatically satisfied.



Duality

A convex optimization problem has a standard form

$$\underset{\mathbf{x}}{\text{minimize}} \quad f_o(\mathbf{x})$$

$$\text{subject to} \quad f_i(\mathbf{x}) \leq 0, \quad i=1, \dots, m$$

$$h_i(\mathbf{x}) = 0, \quad i=1, \dots, p$$

f : convex

h : affine

Lagrangian is defined as

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f_o(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p \nu_i h_i(\mathbf{x})$$

Lagrange dual function

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu})$$

$$= \inf_{\mathbf{x}} \left(f_o(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p \nu_i h_i(\mathbf{x}) \right)$$

Dual problem

$$\underset{\boldsymbol{\lambda}, \boldsymbol{\nu}}{\text{maximize}} \quad g(\boldsymbol{\lambda}, \boldsymbol{\nu})$$

subject to $\boldsymbol{\lambda} \in \mathbb{R}^m$

subject to $\mathbf{A} \mathbf{x} = \mathbf{b}$

Weak duality (always holds)

$$d^* < p^*$$

where d^* is the dual optimum and p^* is the primal optimum.

$$\begin{aligned} \forall \lambda, \nu, \quad g(\lambda, \nu) &= \inf_{\mathbf{x}} \left(f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p \nu_i h_i(\mathbf{x}) \right) \\ &\leq f_0(\mathbf{x}^*) + \underbrace{\sum_{i=1}^m \lambda_i f_i(\mathbf{x}^*)}_{\geq 0} + \underbrace{\sum_{i=1}^p \nu_i h_i(\mathbf{x}^*)}_{\leq 0} \quad [\text{because LHS is infimum}] \\ &\leq f_0(\mathbf{x}^*) \end{aligned}$$

Strong duality (i.e., $d^* = p^*$)

- Slater's condition (convex + strictly feasible) \Rightarrow strong duality
- **(KKT necessary conditions)** For a differentiable objective, if strong duality holds with $\mathbf{x}^*, \lambda^*, \nu^*$ being the primal and dual optimal points, then

$$\begin{array}{lll} f_i(\mathbf{x}^*) \leq 0 & i=1, \dots, m & \text{primal feasibility} \\ h_i(\mathbf{x}^*) = 0 & i=1, \dots, p & \\ \lambda_i^* \geq 0 & i=1, \dots, m & \text{dual feasibility} \\ \lambda_i^* f_i(\mathbf{x}^*) = 0 & i=1, \dots, m & \text{KKT complementary condition} \\ \nabla_{\mathbf{x}} L(\mathbf{x}, \lambda^*, \nu^*) \Big|_{\mathbf{x}=\mathbf{x}^*} = 0 & & \text{Lagrangian gradient vanishing} \end{array}$$

- **(KKT sufficiency for convex problems)** If the primal is convex, then KKT is sufficient. That is, for any $\tilde{\mathbf{x}}, \tilde{\lambda}, \tilde{\nu}$ satisfying KKT conditions, $\tilde{\mathbf{x}}$ is the primal optimum; $\tilde{\lambda}$ and $\tilde{\nu}$ are dual optimum, and strong duality holds.

$$\begin{aligned} g(\tilde{\lambda}, \tilde{\nu}) &= L(\tilde{\mathbf{x}}, \tilde{\lambda}, \tilde{\nu}) \quad [\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda, \nu) \text{ vanishes at } \mathbf{x} = \tilde{\mathbf{x}}] \\ &= f_0(\tilde{\mathbf{x}}) + \sum_{i=1}^m \tilde{\lambda}_i f_i(\tilde{\mathbf{x}}) + \sum_{i=1}^p \tilde{\nu}_i h_i(\tilde{\mathbf{x}}) \\ &= f_0(\tilde{\mathbf{x}}) \end{aligned}$$

Solving the optimization problem for SVM

$$\begin{aligned} \text{Primal:} \quad & \underset{\mathbf{w}, b}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{subject to} \quad 1 - t^{(m)} (\mathbf{w}^T \mathbf{x}^{(m)} + b) \leq 0 \quad \text{for } m=1, \dots, M \end{aligned}$$

$$\text{Lagrangian:} \quad L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{m=1}^M \alpha^{(m)} [1 - t^{(m)} (\mathbf{w}^T \mathbf{x}^{(m)} + b)]$$

$$\text{KKT gradient vanishing} \quad \frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) \stackrel{\text{set}}{=} \mathbf{0} \quad (1)$$

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \alpha) = 0 \quad (2)$$

(1) implies $\mathbf{w} - \sum_{m=1}^M \alpha^{(m)} t^{(m)} \mathbf{x}^{(m)} = 0$

$$\mathbf{w} = \sum_{m=1}^M \alpha^{(m)} t^{(m)} \mathbf{x}^{(m)} \quad (3)$$

(2) implies $\sum_{m=1}^M \alpha^{(m)} t^{(m)} = 0 \quad (4)$

Put (3) and (4) to the Lagrangian:

$$\begin{aligned} L(\mathbf{w}, b, \alpha) &= \frac{1}{2} \underbrace{\|\mathbf{w}\|^2}_{\mathbf{w}^\top \mathbf{w}} + \sum_{m=1}^M \alpha^{(m)} [1 - t^{(m)} (\mathbf{w}^\top \mathbf{x}^{(m)} + b)] \\ &= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M \alpha^{(n)} \alpha^{(m)} t^{(n)} t^{(m)} \mathbf{x}^{(n)}^\top \mathbf{x}^{(m)} \\ &\quad + \sum_{m=1}^M \alpha^{(m)} - \sum_{m=1}^M \alpha^{(m)} t^{(m)} \sum_{n=1}^N \alpha^{(n)} t^{(n)} [\mathbf{x}^{(n)}]^\top \mathbf{x}^{(m)} - b \underbrace{\sum_{m=1}^M \alpha^{(m)} t^{(m)}}_{=0 \text{ by (4)}} \\ &= \sum_{m=1}^M \alpha^{(m)} - \sum_{n=1}^N \sum_{m=1}^M \alpha^{(n)} \alpha^{(m)} t^{(n)} t^{(m)} [\mathbf{x}^{(n)}]^\top \mathbf{x}^{(m)} \end{aligned}$$

Dual:

$$\text{maximize } \sum_{m=1}^M \alpha^{(m)} - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M \alpha^{(n)} \alpha^{(m)} t^{(n)} t^{(m)} [\mathbf{x}^{(n)}]^\top \mathbf{x}^{(m)}$$

subject to $\alpha_i \geq 0$

$$\sum_{m=1}^M \alpha^{(m)} t^{(m)} = 0 \quad \left[\begin{array}{l} \text{Due to } \nabla_b L = 0 \\ \text{Otherwise } \min_{\mathbf{w}, b} L(\mathbf{w}, b, \lambda) = -\infty \end{array} \right]$$

Sequential minimal optimization (SMO)

Loop until converging:

Pick $\alpha^{(i)}, \alpha^{(j)}$

Represent $\alpha^{(i)} = -t^{(i)} \sum_{m \neq i} \alpha^{(m)} t^{(m)}$

Optimize $\alpha^{(j)}$ in the dual

Project back to feasible set

Primal optimal point: If $\alpha_x^{(m)}$ is solved

$$\mathbf{w}_* = \sum_{m=1}^M \alpha_x^{(m)} t^{(m)} \mathbf{x}^{(m)}$$

$$b_* = \frac{1}{2} \left[\min_{i: t^{(i)}=1} \mathbf{w}_*^\top \mathbf{x}^{(i)} + \max_{i: t^{(i)}=-1} \mathbf{w}_*^\top \mathbf{x}^{(i)} \right]$$

Decision boundary

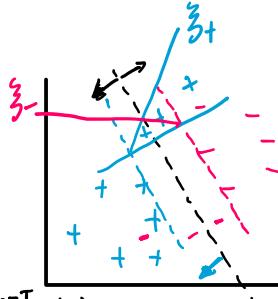
$$\mathbf{w}_*^\top \mathbf{x} + b_* = 0$$

$$\sum_{m=1}^M \alpha_x^{(m)} t^{(m)} [\mathbf{x}^{(m)}]^\top \mathbf{x} + b_* = 0$$

SVM with slack variables

- The above SVM only works for linearly separable problems
- Insensitive to data far away from the margin (dual sparsity)
- Sensitive to data on the margin (support vectors)

$$\begin{aligned}
 & \underset{\mathbf{w}, b}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{m=1}^M \xi^{(m)} \\
 & \text{subject to} \quad t^{(m)} (\mathbf{w}^\top \mathbf{x}^{(m)} + b) \geq 1 - \xi^{(m)} \\
 & \quad \xi^{(m)} \geq 0 \\
 & \text{Dual (C-SVM) [HW]} \\
 & \underset{\alpha}{\text{maximize}} \quad \sum_{m=1}^M \alpha^{(m)} - \frac{1}{2} \sum_{n=1}^M \sum_{m=1}^M \alpha^{(n)} \alpha^{(m)} t^{(n)} t^{(m)} [\mathbf{x}^{(n)}]^\top \mathbf{x}^{(m)} \\
 & \text{subject to} \quad 0 \leq \alpha^{(m)} \leq C \\
 & \quad \sum_{m=1}^M \alpha^{(m)} t^{(m)} = 0
 \end{aligned}$$



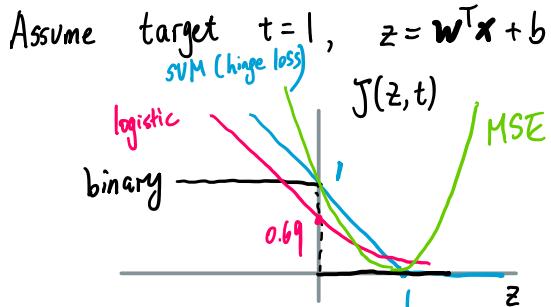
Alternative formulation (ν -SVM)

$$\begin{aligned}
 & \underset{\alpha}{\text{maximize}} \quad -\frac{1}{2} \sum_{n=1}^M \sum_{m=1}^M \alpha^{(n)} \alpha^{(m)} t^{(n)} t^{(m)} [\mathbf{x}^{(n)}]^\top \mathbf{x}^{(m)} \\
 & \text{subject to} \quad 0 \leq \alpha^{(m)} \leq \frac{1}{M} \\
 & \quad \sum_{m=1}^M \alpha^{(m)} t^{(m)} = 0 \\
 & \quad \sum_{m=1}^M \alpha^{(m)} \geq \nu
 \end{aligned}$$

$\nu \in (0, 1)$: upper-bound of fraction of support vectors

Ref: <https://www.csie.ntu.edu.tw/~cjlin/papers/nusvmtutorial.pdf>

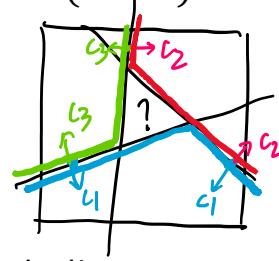
Comparing SVM loss, logistic regression loss, and binary risk



Multi-class SVM

- SVM intrinsically works for binary classification
- Multi-class SVM also works well in practice, but usually with inelegant treatment

- One-vs-one
 - For K -category classification, build $K(K - 1)$ classifiers, and do the voting.
The result may be ambiguous.
- One-vs-the-rest (one-vs-all)
 - Train K classifiers and pick the one category with the largest margin (not normalized by $|\mathbf{w}|$).



Constraints: For any data point $x^{(m)}, t^{(m)}$

$$\begin{aligned} t^{(m)} &= i \in \{1, \dots, K\} \\ \mathbf{w}_i^T \mathbf{x}^{(m)} + b_i - (\mathbf{w}_j^T \mathbf{x}^{(m)} + b_j) &\geq 1 - \xi^{(m)} \quad \forall j \neq i, j \in \{1, \dots, K\} \end{aligned}$$

Kernel tricks

- Idea: Both training and prediction in SVM only depend on inner-product. We may specify the inner-product in some non-linearly transformed space without explicit representation.

Hilbert Space

- **Vector space V**
 - Commutativity* $\forall \mathbf{u}, \mathbf{v} \in V, \mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$
 - Associativity* $\forall \mathbf{u}, \mathbf{v}, \mathbf{w} \in V, (\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$
 - Additive identity* $\exists \mathbf{0} \in V, \forall \mathbf{v} \in V, \mathbf{v} + \mathbf{0} = \mathbf{v}$
 - Additive inverse* $\forall \mathbf{v} \in V, \exists \mathbf{w} \in V, \mathbf{v} + \mathbf{w} = \mathbf{0}$
 - Multiplicative identity* $\forall \mathbf{v} \in V, 1\mathbf{v} = \mathbf{v}$
 - Distributive properties* $\forall a, b \in F, \mathbf{u}, \mathbf{v} \in V$
 $a(\mathbf{u} + \mathbf{v}) = a\mathbf{u} + a\mathbf{v}$ and $(a + b)\mathbf{u} = a\mathbf{u} + b\mathbf{u}$

- **Inner-product space** is a vector space, with inner product
 $\forall \mathbf{u}, \mathbf{v} \in V, \langle \mathbf{u}, \mathbf{v} \rangle \in F$

satisfying the following properties

- Positivity* $\forall \mathbf{v} \in V, \langle \mathbf{v}, \mathbf{v} \rangle \geq 0$
- Definiteness* $\langle \mathbf{v}, \mathbf{v} \rangle = 0$ iff $\mathbf{v} = \mathbf{0}$
- Additivity in first slot* $\forall \mathbf{u}, \mathbf{v}, \mathbf{w} \in V$
 $\langle \mathbf{u} + \mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{u}, \mathbf{w} \rangle + \langle \mathbf{v}, \mathbf{w} \rangle$
[also holds for the second slot]
- Homogeneity in first slot* $\forall a \in F, \forall \mathbf{v}, \mathbf{w} \in V$
 $\langle a\mathbf{v}, \mathbf{w} \rangle = a\langle \mathbf{v}, \mathbf{w} \rangle$
[second slot needs conjugate]
- Conjugate symmetry* $\forall \mathbf{v}, \mathbf{w} \in W, \langle \mathbf{v}, \mathbf{w} \rangle = \overline{\langle \mathbf{w}, \mathbf{v} \rangle}$

The norm $\|\cdot\|: V \rightarrow \mathbb{R}$ is defined as

$$\|\nu\| = \sqrt{\langle \nu, \nu \rangle}$$

- **Hilbert space** is a complete inner-product space

E.g. A finite-dimensional vector space is a Hilbert space.

E.g. $P_m(F)$, m th degree polynomial functions with coefficients in F , equipped with the following inner-product

$$\langle p, q \rangle = \int_0^1 p(x) \overline{q(x)} dx$$

forms a Hilbert space.

Mercer Kernel: $K: X \times X \rightarrow R$

satisfying $\iint f(u)K(u, v)f(v)du dv \geq 0$

for all square-integrable function f , i.e., $\int |f(x)|^2 dx < +\infty$

Discrete analog

$$\mathbf{u}^\top K \mathbf{u} \geq 0 \quad [\text{PSD}]$$

Example. $K(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z}$

Then, $\phi(\mathbf{x}) = \mathbf{x}$

Example (Polynomial kernel). $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z} + 1)^Q$

$$K(\mathbf{x}, \mathbf{z}) = (x_1 z_1 + x_2 z_2 + \dots + x_d z_d + 1)^Q$$

Then, ϕ : all polynomial features up to the degree of Q with some coefficients

Example. Gaussian kernel, Radius Basis Function (RBF) kernel

$$K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|^2)$$

Consider a simple case $\mathbf{x}, \mathbf{z} \in \mathbb{R}$, $\gamma = 1$

$$\begin{aligned} K(\mathbf{x}, \mathbf{z}) &= \exp(-(x - z)^2) \\ &= \exp(-x^2) \exp(-z^2) \exp(2xz) \\ &= \exp(-x^2) \exp(-z^2) \sum_{k=0}^{\infty} \frac{(2xz)^k}{k!} \\ &= \exp(-x^2) \exp(-z^2) \sum_{k=0}^{\infty} \frac{2^k x^k z^k}{k!} \\ &= \sum_{k=0}^{\infty} \left(\sqrt{\frac{2^k}{k!}} \exp(-x^2) x^k \right) \left(\sqrt{\frac{2^k}{k!}} \exp(-z^2) z^k \right) \end{aligned}$$

Thus, ϕ could be

$$x \mapsto \left(\sqrt{\frac{2}{1!}} \exp(-x^2) \cdot x, \sqrt{\frac{2^2}{2!}} \exp(-x^2) \cdot x^2, \dots \right)$$

Note: the feature mapping is not unique.

Reproducing kernel Hilbert space is a set of $f: X \rightarrow \mathbb{R}$

$$\mathcal{H} = \text{span} \{ K(x, \cdot) \mid x \in X \}$$

Reproducing property:

$$\forall f \in \mathcal{H}, \quad \langle f(\cdot), K(\cdot, x) \rangle_{\mathcal{H}} = f(x)$$

$\langle f, \phi_x \rangle_{\mathcal{H}}$

In particular

$$\langle K(\cdot, x), K(\cdot, y) \rangle_{\mathcal{H}} = K(x, y)$$

Operations preserving the kernel

$$k(x, z) = c k_1(x, z) \quad \forall \text{ kernel } k_1, \forall c > 0$$

$$k(x, z) = f(x) K_1(x, z) f(z) \quad \forall f$$

$$K(x, z) = q(k_1(x, z)) \quad \forall \text{ polynomials of nonnegative coefficients}$$

$$K(x, z) = \exp(k_1(x, z))$$

$$K(x, z) = k_1(x, z) + k_2(x, z) \quad \forall \text{ kernel } k_2$$

...

Summary

- Support vectors machines originally work for binary, linearly separable classification.
- Soft SVM is still a linear classifier, but deals with non-linearly separable tasks by slack variables.
- Kernel SVM could be a non-linear classifier for input, which accomplishes linear classification in an implicit Hilbert space.
 - The linear classification is oftentimes not representable by w, b in the Hilbert space, but can be fully specified by support vectors, which could in turn be thought of as model parameters. Thus, the number of parameters of kernel SVM grows with data samples. This is known as a **non-parametric** model.
 - The model capacity shall not be viewed as the dimension of the Hilbert space.
A better measure of SVM's model capacity is the number of support vectors, which is usually small.
[Recall the ν -SVM]

- This gives a clever tradeoff between overfitting (too powerful hypothesis class) and underfitting (too restricted hypothesis class).