

Supervised Learning :: CHEAT SHEET

Linear hard-margin SVM

For labeled data $\mathcal{D} = ((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$, with $y^{(i)} \in \{-1, +1\}$:

- Assume linear separation by $f(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x} + \theta_0$, such that all $+$ -observations are in the positive halfspace $\{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) > 0\}$ and all $-$ -observations are in the negative halfspace $\{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) < 0\}$.
- For a linear separating hyperplane, we have

$$\underbrace{y^{(i)} (\boldsymbol{\theta}^\top \mathbf{x}^{(i)} + \theta_0)}_{=f(\mathbf{x}^{(i)})} > 0 \quad \forall i \in \{1, 2, \dots, n\}.$$

- computes the (signed) distance to the separating hyperplane $f(\mathbf{x}) = 0$, positive for correct classifications, negative for incorrect.

- The distance of f to the whole dataset \mathcal{D} is the smallest distance $\gamma = \min_i \left\{ d(f, \mathbf{x}^{(i)}) \right\}$, which represents the **safety margin**. It is positive if f separates and we want to maximize it.

$$\begin{aligned} \max_{\boldsymbol{\theta}, \theta_0} \quad & \gamma \\ \text{s.t.} \quad & d(f, \mathbf{x}^{(i)}) \geq \gamma \quad \forall i \in \{1, \dots, n\}. \end{aligned}$$

Primal linear hard-margin SVM:

$$\begin{aligned} \min_{\boldsymbol{\theta}, \theta_0} \quad & \frac{1}{2} \|\boldsymbol{\theta}\|^2 \\ \text{s.t.} \quad & y^{(i)} (\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \rangle + \theta_0) \geq 1 \quad \forall i \in \{1, \dots, n\} \end{aligned}$$

This is a convex quadratic program.

Support vectors: All instances $(\mathbf{x}^{(i)}, y^{(i)})$ with minimal margin $y^{(i)} f(\mathbf{x}^{(i)}) = 1$, fulfilling the inequality constraints with equality. All have distance of $\gamma = 1/\|\boldsymbol{\theta}\|$ from the separating hyperplane.

The Lagrange function of the SVM optimization problem is

$$\begin{aligned} L(\boldsymbol{\theta}, \theta_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\boldsymbol{\theta}\|^2 - \sum_{i=1}^n \alpha_i [y^{(i)} (\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \rangle + \theta_0) - 1] \\ \text{s.t.} \quad \alpha_i \geq 0 \quad \forall i \in \{1, \dots, n\}. \end{aligned}$$

The **dual** form of this problem is $\max_{\boldsymbol{\alpha}} \min_{\boldsymbol{\theta}, \theta_0} L(\boldsymbol{\theta}, \theta_0, \boldsymbol{\alpha})$.

We find the stationary point of $L(\boldsymbol{\theta}, \theta_0, \boldsymbol{\alpha})$ w.r.t. $\boldsymbol{\theta}, \theta_0$ and obtain

$$\boldsymbol{\theta} = \sum_{i=1}^n \alpha_i y^{(i)} \mathbf{x}^{(i)}, 0 = \sum_{i=1}^n \alpha_i y^{(i)} \quad \forall i \in \{1, \dots, n\}.$$

Dual linear hard-margin SVM:

$$\begin{aligned} \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y^{(i)} = 0, \\ & \alpha_i \geq 0 \quad \forall i \in \{1, \dots, n\}, \end{aligned}$$

In matrix notation with $\mathbf{K} := \mathbf{X}\mathbf{X}^\top$:

$$\begin{aligned} \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad & \mathbf{1}^\top \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^\top \text{diag}(\mathbf{y}) \mathbf{K} \text{diag}(\mathbf{y}) \boldsymbol{\alpha} \\ \text{s.t.} \quad & \boldsymbol{\alpha}^\top \mathbf{y} = 0, \\ & \boldsymbol{\alpha} \geq 0, \end{aligned}$$

Solution (if existing):

$$\hat{\boldsymbol{\theta}} = \sum_{i=1}^n \hat{\alpha}_i y^{(i)} \mathbf{x}^{(i)}, \quad \theta_0 = y^{(i)} - \langle \hat{\boldsymbol{\theta}}, \mathbf{x}^{(i)} \rangle.$$

Linear Soft-Margin SVM

Allow violations of the margin constraints via slack vars $\zeta^{(i)} \geq 0$

$$y^{(i)} (\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \rangle + \theta_0) \geq 1 - \zeta^{(i)}$$

Now we have two distinct and contradictory goals:

- Maximize the margin.
- Minimize margin violations.

Primal linear soft-margin SVM:

$$\begin{aligned} \min_{\boldsymbol{\theta}, \theta_0, \zeta^{(i)}} \quad & \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{i=1}^n \zeta^{(i)} \\ \text{s.t.} \quad & y^{(i)} (\langle \boldsymbol{\theta}, \mathbf{x}^{(i)} \rangle + \theta_0) \geq 1 - \zeta^{(i)} \quad \forall i \in \{1, \dots, n\}, \\ & \text{and } \zeta^{(i)} \geq 0 \quad \forall i \in \{1, \dots, n\}, \end{aligned}$$

where the constant $C > 0$ controls trade-off between the two conflicting objectives of maximizing the size of the margin and minimizing the frequency and size of margin violations.

Dual linear soft-margin SVM:

$$\begin{aligned} \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \forall i \in \{1, \dots, n\} \quad \text{and} \quad \sum_{i=1}^n \alpha_i y^{(i)} = 0 \end{aligned}$$

- Non-SVs have $\alpha_i = 0$ ($\Rightarrow \mu_i = C \Rightarrow \zeta^{(i)} = 0$) and can be removed from the problem without changing the solution. Their margin $y f(\mathbf{x}) \geq 1$. They are always classified correctly and are never inside of the margin.
- SVs with $0 < \alpha_i < C$ ($\Rightarrow \mu_i > 0 \Rightarrow \zeta^{(i)} = 0$) are located exactly on the margin and have $y f(\mathbf{x}) = 1$.
- SVs with $\alpha_i = C$ have an associated slack $\zeta^{(i)} \geq 0$. They can be on the margin or can be margin violators with $y f(\mathbf{x}) < 1$ (they can even be misclassified if $\zeta^{(i)} \geq 1$).

Regularized ERM representation with hinge loss:

$$\mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) = \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{i=1}^n L(y^{(i)}, f(\mathbf{x}^{(i)})); \quad L(y, f(\mathbf{x})) = \max(1 - y f(\mathbf{x}), 0)$$

Optimization

Algorithm 1 Stochastic subgradient descent (without intercept θ_0)

- 1: **for** $t = 1, 2, \dots$ **do**
- 2: Pick step size α
- 3: Randomly pick an index i
- 4: If $y^{(i)} f(\mathbf{x}^{(i)}) < 1$ set $\boldsymbol{\theta}^{[t+1]} = (1 - \lambda \alpha) \boldsymbol{\theta}^{[t]} + \alpha y^{(i)} \mathbf{x}^{(i)}$
- 5: If $y^{(i)} f(\mathbf{x}^{(i)}) \geq 1$ set $\boldsymbol{\theta}^{[t+1]} = (1 - \lambda \alpha) \boldsymbol{\theta}^{[t]}$
- 6: **end for**

Algorithm 2 Pairwise coordinate ascent in the dual

- 1: Initialize $\boldsymbol{\alpha} = 0$ (or more cleverly)
- 2: **for** $t = 1, 2, \dots$ **do**
- 3: Select some pair α_i, α_j to update next
- 4: Optimize dual w.r.t. α_i, α_j , while holding α_k ($k \neq i, j$) fixed
- 5: **end for**

Supervised Learning :: CHEAT SHEET

Kernel

Kernel = Feature Map + Inner product

Mercer Kernel

A **(Mercer) kernel** on a space \mathcal{X} is a continuous function

$$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

of two arguments with the properties

- Symmetry: $k(\mathbf{x}, \tilde{\mathbf{x}}) = k(\tilde{\mathbf{x}}, \mathbf{x})$ for all $\mathbf{x}, \tilde{\mathbf{x}} \in \mathcal{X}$.
- Positive definiteness: For each finite subset $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ the **kernel Gram matrix** $\mathbf{K} \in \mathbb{R}^{n \times n}$ with entries $K_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is positive semi-definite.

Reproducing property: for all kernels, there must exist a Hilbert space, where a map ϕ of this space satisfies $k(\mathbf{x}, \tilde{\mathbf{x}}) = \langle \phi(\mathbf{x}), \phi(\tilde{\mathbf{x}}) \rangle$. The space is called **reproducing kernel Hilbert space** (RKHS).

Typical Kernels

A kernel can be constructed from other kernels k_1 and k_2 :

- For $\lambda \geq 0$, $\lambda \cdot k_1$ is a kernel.
- $k_1 + k_2$ is a kernel.
- $k_1 \cdot k_2$ is a kernel (thus also k_1^n).

Useful kernels:

- Every constant function taking a non-negative value.
- **Linear kernel:** $k(\mathbf{x}, \tilde{\mathbf{x}}) = \mathbf{x}^\top \tilde{\mathbf{x}}$.
- **Polynomial kernel:** $k(\mathbf{x}, \tilde{\mathbf{x}}) = (\mathbf{x}^\top \tilde{\mathbf{x}} + b)^d$, for $b \geq 0, d \in \mathbb{N}$.

$$\phi(\mathbf{x}) = \left(\sqrt{\binom{d}{k_1, \dots, k_{p+1}}} x_1^{k_1} \dots x_p^{k_p} b^{k_{p+1}/2} \right)_{k_i \geq 0, \sum_i k_i = d}$$

- **Gaussian kernel:** $k(\mathbf{x}, \tilde{\mathbf{x}}) = \exp\left(-\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|^2}{2\sigma^2}\right)$ or $k(\mathbf{x}, \tilde{\mathbf{x}}) = \exp(-\gamma \|\mathbf{x} - \tilde{\mathbf{x}}\|^2)$, $\gamma > 0$

Dual kernelized soft-margin SVM:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \forall i \in \{1, \dots, n\} \quad \text{and} \quad \sum_{i=1}^n \alpha_i y^{(i)} = 0 \end{aligned}$$

Kernel representation of separating hyperplane:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y^{(i)} k(\mathbf{x}^{(i)}, \mathbf{x}) + \theta_0$$

Hyperparameters of SVM

SVMs are somewhat sensitive to its hyperparameters and should always be tuned.

- The choice of C, the choice of the kernel, the kernel parameters are all up to the user.
- Small C allows for margin-violating points in favor of a large margin.
- Large C penalizes margin violators, decision boundary is more wiggly.