

Problem 7: What is the connection between soft-margin SVM and logistic regression?

SVM

$$\min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^N z_i$$

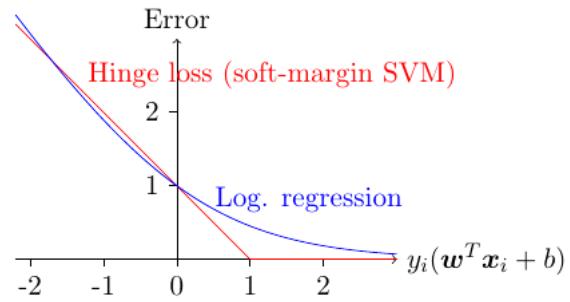
$$y_i (w^T x_i + b) - 1 + z_i \geq 0$$

$$z_i \geq 0$$

Hinge loss
formulation

$$\min_{w, b} E(w, b, C) = \frac{1}{2C} w^T w + \sum_{i=1}^N \text{hinge}(y_i (w^T x_i + b))$$

loss



$$\text{hinge}(z) = \max(0, 1-z)$$

Logistic Regression:

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$p(y_i=1 | x_i, (w, b)) = \sigma(w^T x_i + b)$$

$$p(y_i=-1 | x_i, (w, b)) = 1 - \sigma(w^T x_i + b) = \sigma(-(w^T x_i + b))$$

$$\rightarrow p(y_i | x_i, (w, b)) = \sigma(y_i (w^T x_i + b))$$

Altogether, we get $p(y | X, (w, b)) = \prod_{i=1}^N \sigma(y_i (w^T x_i + b))$

$$\begin{aligned} \min_{w, b} E(w, b) &= -\ln(p(y | X, (w, b))) = -\sum_{i=1}^N \ln((1 + e^{-y_i (w^T x_i + b)})^{-1}) \\ &= \sum_{i=1}^N \ln(1 + e^{-y_i (w^T x_i + b)}) \end{aligned}$$

add ℓ_2 -regularizer

$$\min_{w, b} E(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \ln(1 + e^{-y_i (w^T x_i + b)})$$

loss

Problem 8: Consider a soft-margin SVM fitted to a linearly separable dataset \mathcal{D} using the Hinge loss formulation of the optimization task.

$$\text{minimize}_{\mathbf{w}, b} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \max\{0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)\}$$

- a) Is it guaranteed that all training samples in \mathcal{D} will be assigned the correct label by the model?
- b) Prove that if for some $C_0 \geq 0$ the resulting model classifies all training samples correctly then it will also be the case if we train the model with any larger $C > C_0$.

b.) Denote by

$$h(\mathbf{w}, b) = \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

and let $C > C_0 \geq 0$. Note that

all points lie outside of the margin if and only if $h(\mathbf{w}, b) = 0$.

! and all samples lie outside of the margin

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C h(\mathbf{w}, b)$$

$$(\mathbf{v}^*, d^*) = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C h(\mathbf{w}, b)$$

$$\Rightarrow \frac{1}{2} \|\mathbf{w}^*\|^2 + C h(\mathbf{w}^*, b^*) \leq \frac{1}{2} \|\mathbf{v}^*\|^2 + C h(\mathbf{v}^*, d^*)$$

$$\frac{1}{2} \|\mathbf{v}^*\|^2 + C h(\mathbf{v}^*, d^*) \leq \frac{1}{2} \|\mathbf{w}^*\|^2 + C h(\mathbf{w}^*, b^*)$$

Sum the two inequalities up

$$\Rightarrow \cancel{\frac{1}{2} \|\mathbf{w}^*\|^2} + \cancel{\frac{1}{2} \|\mathbf{v}^*\|^2} + C h(\mathbf{w}^*, b^*) + C h(\mathbf{v}^*, d^*) \leq$$

$$\cancel{\frac{1}{2} \|\mathbf{w}^*\|^2} + \cancel{\frac{1}{2} \|\mathbf{v}^*\|^2} + C h(\mathbf{v}^*, d^*) + C h(\mathbf{w}^*, b^*)$$

$$\Rightarrow (C - C_0) (h(\mathbf{v}^*, d^*) - h(\mathbf{w}^*, b^*)) \leq 0$$

$$\begin{aligned} & \text{---} \\ & \Rightarrow h(v^*, d^*) \leq h(w^*, b^*) = 0 \\ & \Rightarrow h(v^*, d^*) = 0 \quad \Rightarrow 1 - y_i(v^{*\top} x_i + d^*) \leq 0 \end{aligned}$$

Problem 10: Consider the Gaussian kernel

$$\phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2)$$

$$k_G(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}\right), \text{ with } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d.$$

$$\phi(\theta(\mathbf{x}))$$

- a) Suppose you have found a feature map $\theta : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1}$ that transforms your data into a feature space in which a SVM with a Gaussian kernel works well. However, computing $\theta(\mathbf{x})$ is computationally expensive and luckily you discover an efficient method to compute the scalar product

$$k(\mathbf{x}_1, \mathbf{x}_2) = \theta(\mathbf{x}_1)^T \theta(\mathbf{x}_2)$$

$$\mathbf{x} \xrightarrow{\theta} \phi(\mathbf{x})$$

in your feature space without having to compute $\theta(\mathbf{x}_1)$ and $\theta(\mathbf{x}_2)$ explicitly. Show how you can use the scalar product $k(\mathbf{x}_1, \mathbf{x}_2)$ to efficiently compute the Gaussian kernel $k_G(\theta(\mathbf{x}_1), \theta(\mathbf{x}_2))$ in your feature space.

$$k_G(\phi(\mathbf{x}_1), \phi(\mathbf{x}_2)) = \exp\left(-\frac{\|\phi(\mathbf{x}_1) - \phi(\mathbf{x}_2)\|^2}{2\sigma^2}\right)$$

$$\|\mathbf{a} - \mathbf{b}\|^2 = (\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b}) = \mathbf{a}^T \mathbf{a} - 2\mathbf{a}^T \mathbf{b} + \mathbf{b}^T \mathbf{b}$$

$$= \exp\left(-\frac{\phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) - 2\phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) + \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_2)}{2\sigma^2}\right) =$$

$$= \exp\left(-\frac{k(\mathbf{x}_1, \mathbf{x}_1) - 2k(\mathbf{x}_1, \mathbf{x}_2) + k(\mathbf{x}_2, \mathbf{x}_2)}{2\sigma^2}\right)$$

Problem 10: Consider the Gaussian kernel

$$k_G(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}\right), \text{ with } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d.$$

- b) One of the nice things about kernels is that new kernels can be constructed out of already given ones. Use the five kernel construction rules from the lecture to prove that k_G is a kernel.

Hint: Use the Taylor expansion of the exponential function to prove that $\exp \circ k_1$ is a kernel if k_1 is a kernel. Also, consider $k_2(\phi(\mathbf{x}_1), \phi(\mathbf{x}_2))$ with the linear kernel $k_2(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{x}_2$ and a feature map ϕ with only one feature.

Let $k_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be kernels, with $\mathcal{X} \subseteq \mathbb{R}^N$.

Then the following functions k are kernels as well:

- $k(\mathbf{x}_1, \mathbf{x}_2) = k_1(\mathbf{x}_1, \mathbf{x}_2) + k_2(\mathbf{x}_1, \mathbf{x}_2)$ ①
- $k(\mathbf{x}_1, \mathbf{x}_2) = c \cdot k_1(\mathbf{x}_1, \mathbf{x}_2)$, with $c > 0$ ②
- $k(\mathbf{x}_1, \mathbf{x}_2) = k_1(\mathbf{x}_1, \mathbf{x}_2) \cdot k_2(\mathbf{x}_1, \mathbf{x}_2)$ ③
- $k(\mathbf{x}_1, \mathbf{x}_2) = k_3(\phi(\mathbf{x}_1), \phi(\mathbf{x}_2))$, with the kernel k_3 on $\mathcal{X}' \subseteq \mathbb{R}^M$ and $\phi : \mathcal{X} \rightarrow \mathcal{X}'$ ④
- $k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{A} \mathbf{x}_2$, with $\mathbf{A} \in \mathbb{R}^N \times \mathbb{R}^N$ symmetric and positive semidefinite ⑤
- $k_i(x_1, x_2)$ are valid kernels and $k(x_1, x_2) = \lim_{i \rightarrow \infty} k_i(x_1, x_2)$ is also a valid kernel. ⑥

First we show that $\exp(k(x_1, x_2))$ is a kernel if $k(x_1, x_2)$ is a kernel.

$$\exp(k(x_1, x_2)) = 1 + \sum_{n=1}^{\infty} \frac{1}{n!} k(x_1, x_2)^n = \\ = 1 + \lim_{N \rightarrow \infty} \sum_{n=1}^N \frac{1}{n!} k(x_1, x_2)^n$$

$$\textcircled{3} \rightarrow k(x_1, x_2) \cdot k(x_1, x_2) = k(x_1, x_2)^2 \text{ is a kernel} \\ k(x_1, x_2)^n \text{ is a kernel}$$

$$\textcircled{2} \rightarrow \frac{1}{n!} k(x_1, x_2)^n \text{ is a kernel} \\ \textcircled{1} \rightarrow \sum_{n=1}^{\infty} \frac{1}{n!} k(x_1, x_2)^n \text{ is a kernel}$$

$$\textcircled{1} \rightarrow \sum_{n=1}^{\infty} \frac{1}{n!} k(x_1, x_2)^n \text{ is a kernel}$$

$$\textcircled{2} \rightarrow \sum_{n=1}^{\infty} \frac{1}{n!} k(x_1, x_2)^n = \dots$$

$$\textcircled{3} \rightarrow 1 + \sum_{n=1}^{\infty} \frac{1}{n!} k(x_1, x_2)^n = \exp(k(x_1, x_2)) \text{ is a kernel} \checkmark$$

\textcircled{4}) \rightarrow k_a(x_1, x_2) = \exp\left(-\frac{x_1^T x_1}{2\sigma^2}\right) \exp\left(-\frac{x_2^T x_2}{2\sigma^2}\right) \exp\left(\frac{x_1^T x_2}{\sigma^2}\right)

\textcircled{2}, \textcircled{5}) \rightarrow k(x_1, x_2) = x_1^T x_2 \frac{1}{\sigma^2} \text{ is a kernel}

$$\Rightarrow \exp\left(\frac{1}{\sigma^2} x_1^T x_2\right) \text{ is a kernel}$$

\textcircled{4}) \rightarrow k(x_1, x_2) = k_3(\phi(x_1), \phi(x_2)) \text{ is a kernel, take}

$$k_3(x_1, x_2) = x_1 \cdot x_2 \text{, } \phi(x) = \exp\left(-\frac{x^T x}{2\sigma^2}\right)$$

\textcircled{3}) \rightarrow k_a(x_1, x_2) \text{ is a kernel}

Problem 10: Consider the Gaussian kernel

$$k_G(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}\right), \text{ with } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d.$$

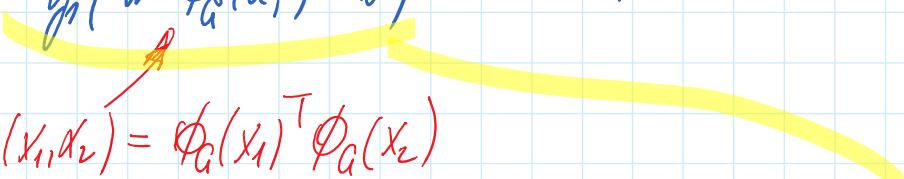
- c) Can any finite set of points be linearly separated in the feature space of the Gaussian kernel if σ can be chosen freely?

Consider the limit $\sigma \rightarrow 0$.

$$k_G(\mathbf{x}_1, \mathbf{x}_2; \sigma) \xrightarrow{\sigma \rightarrow 0} k(\mathbf{x}_1, \mathbf{x}_2) = \begin{cases} 1 & \text{if } \mathbf{x}_1 = \mathbf{x}_2 \\ 0 & \text{if } \mathbf{x}_1 \neq \mathbf{x}_2 \end{cases}$$

All training samples are correctly classified if

$$y_i (\mathbf{w}^T \phi_G(\mathbf{x}_i) + b) > 0 \quad \forall i$$



$$k_G(\mathbf{x}_1, \mathbf{x}_2) = \phi_G(\mathbf{x}_1)^T \phi_G(\mathbf{x}_2)$$

Substitute $\mathbf{w} = \sum_j y_j \alpha_j \phi_G(\mathbf{x}_j)$ into



$$y_i \left(\underbrace{\sum_j y_j \alpha_j \phi_G(\mathbf{x}_j)^T \phi_G(\mathbf{x}_i)}_{k_G(\mathbf{x}_j, \mathbf{x}_i; \sigma)} + b \right) > 0$$



$$\xrightarrow{\sigma \rightarrow 0} \underbrace{y_i^2}_{1} \alpha_i + y_i b = \alpha_i + y_i b$$

By choosing $b=0$ and $\alpha_i > 0$ we can fulfill
 $\underline{\sigma \text{ very small}}$

