

kernel classification and regression

- ridge regression

given data and response $(x_i, y_i)_{i=1}^n$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}_{n \times D} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

linear regression $x_i^T \beta \approx y_i$, $\beta \in \mathbb{R}^D$

regularize by $\|\beta\|_2^2$

$$\hat{\beta} := \underset{\beta \in \mathbb{R}^D}{\operatorname{argmin}} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_2^2, \lambda > 0$$

$$L(\beta) = \|X\beta - y\|_2^2 + \lambda \|\beta\|_2^2 \\ = \beta^T X^T X \beta - 2y^T X \beta + \|y\|_2^2 + \lambda \|\beta\|_2^2$$

$$\frac{\partial L}{\partial \beta} = 2((X^T X) \beta - X^T y + \lambda \beta) = 0$$

$$(X^T X + \lambda I) \hat{\beta} = X^T y$$

$$(X^T X)_{D \times D} \quad \text{PSD}$$

$(X^T X + \lambda I)$ invertible

$$\Rightarrow \hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$$

The predictive model: $x^{(te)}$ new data sample

$$\hat{y}(x^{(te)}) = x^{(te) T} \hat{\beta} \\ = x^{(te) T} (X^T X + \lambda I)^{-1} (X^T y) \\ = y^T \underbrace{X (X^T X + \lambda I)^{-1} x^{(te)}}_{n \times n} \\ \stackrel{\text{(claim)}}{=} y^T \underbrace{(X X^T + \lambda I)^{-1} X}_{n \times n} x^{(te)}$$

If suppose $\begin{bmatrix} X \\ \vdots \\ x^{(te)} \end{bmatrix}_{n+1 \times D} = \begin{bmatrix} U \\ \vdots \\ U \end{bmatrix}_{n+1 \times D} \begin{bmatrix} S & V^T \end{bmatrix}_{D \times D}$

$$X^T X = V S^2 V^T \quad D \times D$$

$$X^T X^T = U S^2 U^T \quad n \times n \quad (\text{rank } \leq D)$$

$$\hat{y} = X (X^T X + \lambda I)^{-1} = U S V^T (V (S^2 + \lambda I) V^T)^{-1}$$

$$= U S (S^2 + \lambda I)^{-1} V^T$$

$$\hat{y} = (X^T X + \lambda I)^{-1} X = (U (S^2 + \lambda I) U^T)^{-1} U S V^T$$

$$= U (S^2 + \lambda I)^{-1} S V^T$$

bcz S is diagonal, $\hat{y} = \hat{y}$. $\#$

Thus

$$\hat{y}(x^{(te)}) = y^T \underbrace{(X X^T + \lambda I)^{-1} X}_{n \times n} x^{(te)} \\ = \sum_{i=1}^n ((G + \lambda I)^{-1} g_{te})_i y_i \quad (*)$$

$$G := X X^T, G_{i,j} = x_i^T x_j \quad \text{Gram matrix } n \times n$$

$$g_{te} := X^T x^{(te)}, (g_{te})_i = x_i^T x^{(te)}, i=1, \dots, n$$

$\hat{y}(x^{(te)})$ is a weighted avg of y_i (training labels) and only involves

$$x_i^T x_j, x_i^T x^{(te)}$$

idea: change $x^T y$ to $\langle \phi(x), \phi(y) \rangle_{\mathcal{H}}$

equivalently, use feature map $\phi(x_i)$ in place of x_i in the regression.

If we just change the inner product, we get

$$\hat{y}(x^{(te)}) = y^T (G + \lambda I)^{-1} g_{te}, \quad (*)$$

$$G = (k(x_i, x_j))_{i,j=1}^n$$

$$g_{te} = (k(x_i, x^{(te)}))_{i=1}^n$$

We show that $(*)$ is same as if we use $\phi(x)$ in the regression.

the model is to find $\beta \in \mathcal{H}$,

$$\langle \phi(x_i), \beta \rangle_{\mathcal{H}} \approx y_i, \quad i=1, \dots, n$$

$$\min_{\beta \in \mathcal{H}} \sum_{i=1}^n (\langle \phi(x_i), \beta \rangle_{\mathcal{H}} - y_i)^2 + \lambda \|\beta\|_{\mathcal{H}}^2$$

thus $\beta = \sum_{j=1}^n \alpha_j \phi(x_j)$, $\alpha \in \mathbb{R}^n$ TBD.

$$\langle \phi(x_i), \beta \rangle_{\mathcal{H}} = \sum_{j=1}^n \alpha_j \underbrace{\langle \phi(x_i), \phi(x_j) \rangle_{\mathcal{H}}}_{k(x_i, x_j)} \\ = (\alpha^T G)_{i,i} = \alpha^T g_{te}$$

$$\|\beta\|_{\mathcal{H}}^2 = \langle \sum_{j=1}^n \alpha_j \phi(x_j), \sum_{j=1}^n \alpha_j \phi(x_j) \rangle_{\mathcal{H}} \\ = \alpha^T G \alpha,$$

$$= \alpha^T G \alpha,$$

$$\min_{\alpha} \|\alpha^T G \alpha - y\|_2^2 + \lambda \|\alpha\|_2^2 = L(\alpha)$$

$$\frac{\partial L}{\partial \alpha} = 2(G^T G \alpha - G^T y + \lambda G \alpha) = 0$$

$$\Rightarrow G(G + \lambda I) \alpha = G y, \quad (G = G^T, \text{PSD})$$

$$\Rightarrow \hat{\alpha} = (G + \lambda I)^{-1} y$$

$$\hat{y}(x') = \langle \phi(x'), \hat{\beta} \rangle_{\mathcal{H}} \\ = \sum_{i=1}^n \hat{\alpha}_i \langle \phi(x'), \phi(x_i) \rangle_{\mathcal{H}}$$

$$= \sum_{i=1}^n k(x', x_i) \hat{\alpha}_i = \hat{\alpha}^T g_{te}$$

$$= y^T (G + \lambda I)^{-1} g_{te} \quad \text{which is } (x') \#.$$

- classification

$$\{(x_i, y_i)\}_{i=1}^n, y_i = \pm 1$$

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n l(y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}}^2$$

$$l(y, f(x)) \quad \text{loss eq. logistic loss}$$

$$\Pr[l(x)=1 | x] = \frac{e^{f(x)}}{1+e^{f(x)}}$$

kernel logistic regression

can generalize to M-class classification

multi-logit regression:

$$\Pr[l(x)=m | x] = \frac{e^{f_m(x)}}{\sum_{j=1}^M e^{f_j(x)}},$$

$$f_1, \dots, f_M \in \mathcal{H}$$

This is same as in softmax layer in NN, where f_1, \dots, f_M are last hidden layer activations in an NN.

Kernel logistic regression is related to SVM

- linear SVM

$$y_i = \pm 1$$

$$\min_{a \in \mathbb{R}^D, b \in \mathbb{R}} \|a\|_2^2$$

$$\text{s.t. } y_i(a^T x_i - b) \geq 1$$

Convex optimization

$$a^T x_i - b < -1$$

$$a^T x_i - b = 1$$

$$a^T x_i - b > 1$$

$$a^T x_i - b = 1$$

$$a^T x_i - b > 1$$

relax to hinge loss (soft-margin formulation)

when there is no line to separate the two classes perfectly.

$$\min_{a,b} \sum_{i=1}^n l(y_i, a^T x_i - b) + \lambda \|a\|_2^2$$

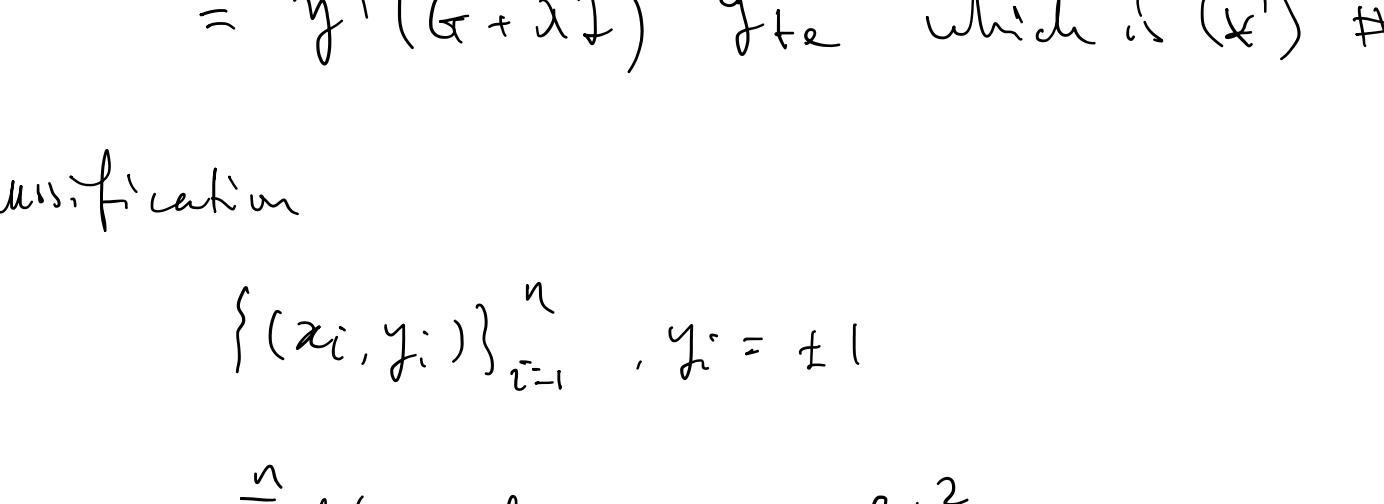
loss + penalty

- kernel SVM

$$\hat{\beta} := \begin{bmatrix} a \\ b \end{bmatrix}, \bar{x}_i := \begin{bmatrix} x_i \\ 1 \end{bmatrix},$$

$\hat{\beta}^T \bar{x}_i \rightarrow$ kernel inner product

$$\hat{\beta}^T \bar{x}_i = \langle \hat{\beta}, \phi(\bar{x}_i) \rangle_{\mathcal{H}}, \hat{\beta} \in \mathcal{H}$$



$$\min_{a,b} \sum_{i=1}^n l(y_i, a^T x_i - b) + \lambda \|a\|_2^2$$

relax to hinge loss (soft-margin formulation)

when there is no line to separate the two classes perfectly.

$$\min_{a,b} \sum_{i=1}^n l(y_i, a^T x_i - b) + \lambda \|a\|_2^2$$

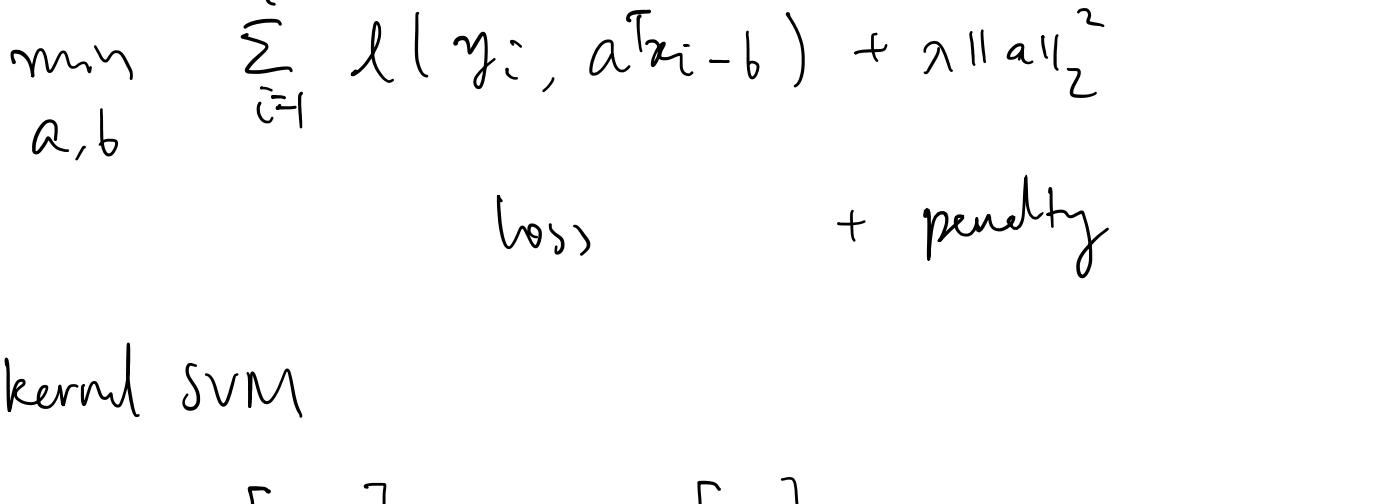
loss + penalty

kernel SVM

$$\hat{\beta} := \begin{bmatrix} a \\ b \end{bmatrix}, \bar{x}_i := \begin{bmatrix} x_i \\ 1 \end{bmatrix},$$

$\hat{\beta}^T \bar{x}_i \rightarrow$ kernel inner product

$$\hat{\beta}^T \bar{x}_i = \langle \hat{\beta}, \phi(\bar{x}_i) \rangle_{\mathcal{H}}, \hat{\beta} \in \mathcal{H}$$



Code implementation in Python sklearn: <https://scikit-learn.org/stable/modules/svm.html>