



Lecture 16

. SVM dual problem . linear C-SVM · Kernel SVM

Next lecture on
 Zoom
 Remote, recorded

Lecture V: Support Vector Machines

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

November, 2022

November, 2022

Linear SVM's

The margin and the expected classification error Maximum Margin Linear classifiers Linear classifiers for non-linearly separable data

Non linear SVM The "kernel trick" Kernels Prediction with SVM

Extensions

Marina Meila: Lecture V: Support Vector Machines

L₁ SVM Multi-class and One class SVM SV Regression

Reading HTF Ch.: Ch. 12.1–3, Murphy Ch.: Ch 14 (14.1,14.2–14.2.4 kernels, 14.4 and equations (14.28,14.29) kernel trick, 14.5.1.–3 Support Vector Machines), Bach Ch.: 7.1–7.4, 7.7 Additional Reading: C. Burges - "A tutorial on SVM for pattern recognition" These notes: Appendices (convex optimization) are optional.

Second idea

The Second idea is to formulate (10) as a quadratic optimization problem.

$$\min_{w,b} \frac{1}{2} ||w||^2 \text{ s.t } y^i(w^T x^i + b) \ge 1 \text{ for all } i = 1:n$$
(13)

This is the Linear SVM (primal) optimization problem

- This problem has a strongly convex objective $||w||^2$, and constraints $y^i(w^Tx^i + b)$ linear in (w, b).
- ▶ Hence this is a convex problem, and can be studied with the tools of convex optimization.
 - solved by off-the shelf optimizer. E.g. CVX

2022

November

The distance of a point x to a hyperplane $H_{w,b}$

$$d(x, H_{w,b}) = \frac{|w^{T}x + b|}{||w||}$$
(14)

Intuition: denote

$$\tilde{w} = \frac{w}{||w||}, \ \tilde{b} = \frac{b}{||w||}, \ x' = \tilde{w}^T x.$$
 (15)

Obviously $H_{w,b} = H_{\tilde{w},\tilde{b}}$, and x' is the length of the projection of point x on the direction of w. The distance is measured along the normal through x to H; note that if $x' = -\tilde{b}$ then $x \in H_{w,b}$ and $d(x, H_{w,b}) = 0$; in general, the distance along this line will be $|x' - (-\tilde{b})|$.



and $b = y^i - w^T x^i$ for any *i* with $\alpha_i > 0$.

- Support vector is a data point x^i such that $\alpha_i > 0$.
- According to (17), the final decision boundary is determined by the support vectors (i.e. does not depend explicitly on any data point that is not a support vector).

vember, 2022

²The derivations of these results are in the Appendix

Karush Kuhn Tuckor (at optimum
$$W^*, U^*, V^*$$
)

$$\begin{array}{c}
\frac{\partial L}{\partial W} = 0 = W - \sum_{i=1}^{n} \alpha_i y^i x^i \Rightarrow W^* = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} y^i x^i \\
\frac{\partial L}{\partial W} = 0 = -\sum_{i=1}^{n} \alpha_i y^i x^i \Rightarrow W^* = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} y^i x^i \\
\frac{\partial L}{\partial W} = 0 = -\sum_{i=1}^{n} \alpha_i y^i x^i \Rightarrow W^* = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} y^i x^i \\
\frac{\partial L}{\partial W} = 0 = -\sum_{i=1}^{n} \alpha_i y^i x^i \Rightarrow W^* = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} y^i x^i \\
\frac{\partial L}{\partial W} = 0 = -\sum_{i=1}^{n} \alpha_i y^i x^i \Rightarrow W^* = \sum_{i=1}^{n} \alpha_i y^i y^i x^i \\
\frac{\partial L}{\partial W} = 0 = -\sum_{i=1}^{n} \alpha_i y^i x^i \Rightarrow W^* \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i + \frac{\partial}{\partial W} \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i x^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^{n} \frac{\partial}{\partial x^i} x^i y^i y^i x^i \\
\frac{\partial L}{\partial W} = \sum_{i=1}^$$

$$\begin{split} \|W\|^{\mathcal{X}} &= W^{T}W = \sum_{\substack{i,j=1\\i \neq i}}^{n} \alpha_{i} \alpha_{j} y^{i} y^{j} (x^{i})^{T} x^{j} \\ \sum_{\substack{i=1\\i \neq i}}^{n} (\alpha_{i} y^{i} x^{i}) W^{T} = -w - \\ (\sum_{\substack{i=1\\i \neq i}}^{n} \alpha_{i} y^{i}) b = 0 \\ \vdots \\ \vdots \end{split}$$

$$g \equiv L(W^*, b^*, \alpha) =$$

Finding
$$6^{\times}$$

1) pick i support vector
2) $(w^{\star})^{T}x^{i} + b = y^{i} \Rightarrow b = y^{i} - (w^{\star})^{T}x^{i}$
 $g(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i,j=1}^{N} \alpha_{i} \alpha_{j} y^{i} y^{j} y^{j} x^{j} x^{j} = 1^{T} \alpha - \frac{1}{2} \alpha^{T} \overline{G} \alpha$
 $1 = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} \in \mathbb{R}^{n}$
 $\overline{G} = \begin{bmatrix} y^{i} y^{j} y^{i} x^{j} \end{bmatrix} \xrightarrow{G} = \begin{bmatrix} x \\ y^{i} y^{j} x^{j} \end{bmatrix}$

Dual SVM optimization problem

- Any convex optimization problem has a dual problem. In SVM, it is both illuminating and practical to solve the dual problem.
- The dual to problem (13) is

$$\max_{\alpha_{1:n}} \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i} \alpha_{i} \alpha_{j} y^{i} y^{j} x^{iT} x_{j} \text{ s.t } \alpha_{i} \ge 0 \text{ for all } i \text{ and } \sum_{i} \alpha_{i} y^{i} = 0.$$
(18)

- This is a <u>quadratic</u> problem with *n* variables on a convex domain.
- Dual problem in matrix form
 - Denote $\alpha = [\alpha_i]_{i=1:n}$, $y = [y^i]_{i=1:n}$, $G_{ij} = x^{iT}x_j$, $\overline{G}_{ij} = y^i y^j x^{iT}x_j$, $G = [G_{ij}] \in \mathbb{R}^{n \times n}$, $\overline{G} = [\overline{G}_{ij}] \in \mathbb{R}^{n \times n}$.

$$\max_{\alpha \in \mathbb{R}^n} \mathbf{1}^T \alpha - \frac{1}{2} \alpha^T \bar{G} \alpha \quad \text{s.t } \alpha \succeq 0 \text{ and } y^T \alpha = 0.$$
(19)

- $g(\alpha) = 1^T \alpha \frac{1}{2} \alpha^T \bar{G} \alpha$ is the dual objective function
- G is called the Gram matrix of the data. Note that $\overline{G} = \operatorname{diag} y^{1:n^T} G \operatorname{diag} y^{1:n}$.
- At the dual optimum

• $\alpha_i > 0$ for constraints that are satisfied with equality, i.e. tight • $\alpha_i = 0$ for the slack constraints

- . Linear classifier, separable data /
- ____, any data /
- Non-linear, --- + Kernel tride

Non-linearly separable problems and their duals penalty 20 The C-SVM minimize_{w,b,ξ} $\frac{1}{2}||w||^2 + C\sum \xi_i$ (20)C small -> Var V s.t. $y^{i}(w^{T}x^{i}+b) \geq 1-\xi_{i}$ $\xi_{i} \geq 0$ slack variable C large → empirical In the above, ξ_i are the slack variables. Dual \mathcal{M} $\sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i} \alpha_{i} \alpha_{j} y^{i} y_{j} x^{iT} x_{j}$ maximize (21)s.t. $C \ge \alpha_i \ge 0$ for all $i \implies \text{More SV}$ $\sum \alpha_i y^i = 0$ \Rightarrow two types of SV \leq $a_i < C$ data point x^i is "on the margin" $\Rightarrow y^i(w^T x^i + b) = 1$ (original SV) $\mathbf{P}_{i} \alpha_{i} = C_{i}$ data point x^{i} cannot be classified with margin 1 (margin error) $\Leftrightarrow y^i(w^Tx^i+b) < 1$ ³Lagrangean $L(w, b, \xi, \alpha, \mu) = \frac{1}{2} ||w||^2 + C \sum_i \xi_i - \sum_i \alpha_i [y^i (w^T x^i + b) - 1 + \xi_i] - \sum_i \mu_i \xi_i$ with $\alpha_i > 0, \ \xi_i > 0, \ \mu_i > 0$

13

The ν -SVM

minimize_{w,b,\xi,\rho}
$$\frac{1}{2}||w||^{2} - \nu\rho + \frac{1}{n}\sum_{i}\xi_{i}$$
 (22)
s.t.
$$y^{i}(w^{T}x^{i} + b) \ge \rho - \xi_{i}$$
 (23)
$$\xi_{i} \ge 0$$
 (24)
$$\rho \ge 0$$
 (25)

where $\nu \in [0, 1]$ is a parameter. Dual⁴:

 $-\frac{1}{2}\sum_{i}\alpha_{i}\alpha_{j}y^{i}y^{j}x^{iT}x^{j}$ ⁽²⁶⁾

t.
$$\frac{1}{n} \ge \alpha_i \ge 0$$
 for all i (27)

$$\sum_{i} \alpha_{i} y^{i} = 0 \tag{28}$$

$$\sum_{i} \alpha_{i} \geq \nu \tag{29}$$

Properties If $\rho > 0$ then:

• ν is an upper bound on #margin errors/n (if $\sum_i \alpha_i = \nu$)

 $\operatorname{maximize}_{\alpha}$

s.

- $\blacktriangleright v$ is a lower bound on #(original support vectors + margin errors)/n
- ν -SVM leads to the same w, b as C-SVM with $C = 1/\nu$

$$\label{eq:Lagrangean L} \begin{split} ^{4} \mathsf{Lagrangean } L(\mathbf{w}, b, \xi, \rho, \alpha, \mu, \delta) \ = \ \frac{1}{2} ||\mathbf{w}||^2 - \nu\rho + \frac{1}{n} \sum_i \xi_i - \sum_i \alpha_i [y^i(\mathbf{w}^T \mathbf{x}^i + b) - \rho + \xi_i] - \sum_i \mu_i \xi_i - \delta\rho \\ \text{with } \alpha_i \geq 0, \ \delta \geq 0, \ \mu_i \geq 0 \end{split}$$

A simple error bound

$$L_{01}(f_n) \leq E\left[\frac{\# \text{support vectors of } f_{n+1}}{n+1}\right]$$

(30)

where f_n denotes the SVM trained on a sample of size n. Exercise Use the Homework 6 to prove this result.

Non-linear SVM

How to use linear classifier on data that is not linearly separable? An old trick

1. Map the data $x^{1:n}$ to a higher dimensional space

 $x \to z = \phi(x) \in \mathcal{H}$, with dim $\mathcal{H} >> n$.

2. Construct a linear classifier $w^T z + b$ for the data in \mathcal{H}

In other words, we are implementing the non-linear classifier

$$f(x) = w^{T} \phi(x) + b = w_{1} \phi_{1}(x) + w_{2} \phi_{2}(x) + \ldots + w_{m} \phi_{m}(x) + b$$
(31)

 $\dim W = \dim \varphi(x)$ $\equiv \dim \mathcal{H}$





- How does SVM change?
 Systematic way to get φ?

Non-linear SV problem

Primal problem minimize ¹/₂ ||w||² s.t yⁱ(w^T \phi(xⁱ) + b) - 1 ≥ 0 for all i.
 Dual problem

$$\max_{\alpha_{1:n}} \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i} \alpha_{i} \alpha_{j} y^{i} y_{j} \phi(x^{i})^{T} \phi(x_{j}) \text{ s.t. } \alpha_{i} \geq 0 \text{ for all } i \text{ and } \sum_{i} y^{i} \alpha_{i} = 0 \quad (32)$$

$$i \quad Xerred \quad (32)$$

$$G_{ij} = \phi(x^{i})^{T} \phi(x^{j}) \text{ and } \overline{G} = y^{T} Gy \quad (33)$$

$$G_{ij} \text{ has been redefined in terms of } \phi \quad only \text{ on } (32)$$

$$\max_{\alpha} 1^{T} \alpha - \frac{1}{2} \alpha^{T} \overline{G} \alpha \quad \text{s.t. } \alpha_{i} \geq 0, y^{T} \alpha = 0 \quad (34)$$
Same as (19)!
$$G \in \mathbb{R}^{n \times n}$$

$$(2^{i})^{T} z^{j} \text{ in } O(d)$$

$$\lim_{\alpha} x^{i} \gamma = \varphi(x)^{T} \varphi(x^{j}) \quad \text{in } \mathcal{R} \text{ bilinear}$$

dim w = dim H

Marina Meila: Lecture V: Support Vector Machines

The "Kernel Trick"

- idea The result (34) is the celebrated kernel trick of the SVM literature. We can make the following remarks.
 - 1. The ϕ vectors enter the SVM optimization problem only trough the Gram matrix, thus only as the scalar products $\phi(x^i)^T \phi(x_i)$. We denote by K(x, x') the function

$$K(x, x') = K(x', x) = \phi(x)^T \phi(x')$$
 (35)

K is called the **kernel** function. If *K* can be computed efficiently, then the Gram matrix *G* can also be computed efficiently. This is exactly what one does in practice: we choose ϕ implicitly by choosing a kernel *K*. Hereby we also ensure that *K* can be computed efficiently.

- 2. Once G is obtained, the SVM optimization is independent of the dimension of x and of the dimension of $z = \phi(x)$. The complexity of the SVM optimization depends only on n the number of examples. This means that we can choose a very high dimensional ϕ without any penalty on the optimization cost.
- 3. Classifying a new point x. As we know, the SVM classification rule is

$$f(x) = w^{T}\phi(x) + b = \sum_{i=1}^{n} \alpha_{i} y^{i} \phi(x^{i})^{T} \phi(x) = \sum_{i=1}^{n} \alpha_{i} y^{i} K(x^{i}, x)$$
(36)

Hence, the classification rule is expressed in terms of the support vectors and the kernel only. No operations other than scalar product are performed in the high dimensional space H.

mber, 2022

Kernels

The previous section shows why SVMs are often called **kernel machines**. If we choose a kernel, we have all the benefits of a mapping in high dimensions, without ever carrying on any operations in that high dimensional space. The most usual kernel functions are

 $K(x,x') = (1 + x^T x')^p$ $K(x,x') = \tanh(\sigma x^T x' - \beta)$ $K(x,x') = e^{-\frac{||x-x'||^2}{\sigma^2}}$ the polynomial kernel of degree *p* the "neural network" kernel

the Gaussian or radial basis function (RBF) kernel it's ϕ is $\infty\text{-dimensional}$

Happy Thanks giving

The Mercer condition

- ► How do we verify that a chosen *K* is is a valid kernel, i.e that there exists a ϕ so that $K(x, x') = \phi(x)^T \phi(x')$?
- This property is ensured by a positivity condition known as the Mercer condition.

Mercer condition

Let (\mathcal{X}, μ) be a finite measure space. A symmetric function $K : \mathcal{X} \times \mathcal{X}$, can be written in the form $K(x, x') = \phi(x)^T \phi(x')$ for some $\phi : \mathcal{X} \to \mathcal{H} \subset \mathbb{R}^m$ iff

 $\int_{\mathcal{X}^2} K(x,x')g(x)g(x')d\mu(x)d\mu(x') \ge 0 \quad \text{for all } g \text{ such that } ||g(x)||_{L_2} < \infty$ (37)

In other words, K must be a positive semidefinite operator on L₂.
 If K satisfies the Mercer condition, there is no guarantee that the corresponding φ is unique, or that it is finite-dimensional.

Quadratic kernel



- C-SVM, polynomial degree 2 kernel, n = 200, C = 10000
- The two ellipses show that a constant shift to the data (xⁱ ← xⁱ + v, v ∈ ℝⁿ) can affect non-linear kernel classifiers.

RBF kernel and Support Vectors



Prediction with SVM

- Estimating b
 - For any *i* support vector, $w^T x^i + b = y^i$ because the classification is tight
 - Alternatively, if there are slack variables, $w^T x^i + b = y^i (1 \xi_i)$
 - $\blacktriangleright \text{ Hence, } b = y^i (1 \xi_i) w^T x^i$
 - For non-linear SVM, where w is not known explicitly, $w = \sum_{j} \alpha_{j} y^{j} \phi(x_{j})$. Hence, $b = y^{i} (1 - \xi_{i}) - \sum_{j=1}^{n} \alpha_{j} y^{j} K(x^{i}, x^{j})$ for any *i* support vector

Given new x

$$\hat{y}(x) = \operatorname{sgn}(w^{T}x + b) = \operatorname{sgn}\left(\sum_{i=1}^{n} \alpha_{i} y^{i} \mathcal{K}(x^{i}, x) + b\right).$$
(38)

L1-SVM

▶ If the regularization $||w||^2$, based on l_2 norm, is replaced with the l_1 norm $||w||_1$, we obtain what is known as the Linear L1-SVM

 $\min_{w,b} ||w||_1 + C \sum_i \xi_i \quad \text{s.t } y^i (w^T x^i + b) \ge 1 - \xi_i, \ \xi_i \ge 0 \text{ for all } i = 1: n$ (39)

- The use of the l_1 norm promotes sparsity in the entries of w
- The Non-linear L1-SVM is

$$f(x) = \sum_{i} (\alpha_{i}^{+} + \alpha_{i}^{-}) y^{i} K(x_{i}, x) + b \quad \text{classifier}$$

$$(40)$$

$$\min_{\alpha \pm, b} \qquad \sum_{i} (\alpha_{i}^{+} + \alpha_{i}^{-}) + C \sum_{i} \xi_{i} \quad \text{s.t } y^{i} f(x^{i}) \ge 1 - \xi_{i}, \ \xi_{i}, \alpha_{i}^{\pm} \ge 0 \text{ for all } i = 1 \quad (41)$$

- This formulation enforces α_i⁺ = 0 or α_i⁻ = 0 for all *i*. If we set w_i = α_i⁺ − α_i⁻, we can write f(x) = ∑_i w_iyⁱK(xⁱ, x) + b, a linear classifier in the non-linear features K(xⁱ, x).
 The L1-SVM problems are Linear Programs
- The dual L1-SVM problems are also linear programs
- The L1-SVM is no longer a Maximum Margin classifier

Multi-class and One class SVM

Multiclass SVM

For a problem with K possible classes, we construct K separating hyperplanes $w_r^T x + b_r = 0$.

minimize

$$\frac{1}{2}\sum_{r=1}^{K}||w_{r}||^{2}+\frac{C}{n}\sum_{i,r}\xi_{i,r}$$
(42)

$$w_{y^i}^T x^i + b_{y^i} \ge w_r^T x^i + b_r + 1 - \xi_{i,r}$$
 for all $i = 1 : n, r \neq y^i$ (43)

$$\xi_{i,r} \geq 0$$
 (44)

One-class SVM This SVM finds the "support regions" of the data, by separating the data from the origin by a hyperplane. It's mostly used with the Gaussian kernel, that projects the data on the unit sphere. The formulation below is identical to the ν -SVM where all points have label 1.

> $\frac{1}{2}||w||^2 - \nu\rho + \frac{1}{n}\sum_{i}\xi_i$ minimize (45)

s.t. $w^T x^i + b \ge \rho - \xi_i$ $\xi_i \ge 0$ $\rho \ge 0$ (46)

$$T_i \geq 0$$
 (47)

(48)

SV Regression

mini

s.

The idea is to construct a "tolerance interval" of $\pm \epsilon$ around the regressor f and to penalize data points for being outside this tolerance margin. In words, we try to construct the smoothest function that goes within ϵ of the data points.

mize
$$\frac{1}{2} ||w||^2 + C \sum_i (\xi_i^+ + \xi_i^-)$$
 (49)

t.
$$\epsilon + \xi_i^+ \ge w^T x^i + b - y^i \ge -\epsilon - \xi_i^-$$
 (50)

$$\xi_i^{\pm} \ge 0$$
 (51)

$$p \ge 0$$
 (52)

The above problem is a linear regression, but with the kernel trick we obtain a kernel regressor of the form $f(x) = \sum_i (\alpha_i^- - \alpha_i^+) K(x^i, x) + b$