

Support Vector Machines 1

Review: The Kernel Trick

- A function $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is an inner product kernel if there exists a function $\Phi(\mathbf{x})$ mapping to an inner product space such that

$$k(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle \quad \forall \mathbf{u}, \mathbf{v}.$$

- A machine learning algorithm is said to be *kernelizable* if it is possible to formulate the algorithm such that all training instances \mathbf{x}_i and any test instance \mathbf{x} occur exclusively in inner products of the form $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$, $\langle \mathbf{x}_i, \mathbf{x} \rangle$ or $\langle \mathbf{x}, \mathbf{x} \rangle$.
- Suppose Φ is a feature map associated to an inner product kernel k .
- If we apply a kernelizable algorithm to the training data

$$(\Phi(\mathbf{x}_1), y_1), \dots, (\Phi(\mathbf{x}_n), y_n)$$

then we can formulate the algorithm such that transformed feature vectors only appear via inner products $\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$ with other transformed feature vectors.

- Can implement by evaluating $k(\mathbf{x}, \mathbf{x}')$

which eliminates the need to ever compute $\Phi(\mathbf{x})$ explicitly.

OSM Hyperplane Classifier

- Our goal today is to kernelize the OSM hyperplane (equivalently, regularized ERM with the hinge loss). This is the linear classifier

$$\mathbf{x} \mapsto \text{sign}\{(\mathbf{w}^*)^T \mathbf{x} + b^*\}$$

where $(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*)$ is a solution of

$$\begin{aligned} \min_{\mathbf{w}, b, \boldsymbol{\xi}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i && \text{(OSM)} \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i \in \{1, \dots, n\} \\ & \xi_i \geq 0 \quad \forall i \in \{1, \dots, n\} \end{aligned}$$

- It's not obvious how to kernelize this method.
- To kernelize it, we will apply the theory of constrained optimization.

Constrained Optimization

- A constrained optimization problem has the form

$$\begin{array}{ll} \min_{u \in \mathbb{R}^d} & f(u) \\ \text{s.t.} & g_i(u) \leq 0 \quad \forall i \in \{1, \dots, r\} \\ & h_j(u) = 0 \quad \forall j \in \{1, \dots, s\} \end{array}$$

- If u satisfies all of the constraints, it is said to be *feasible*.

$$\text{feasible set} := \{u : g_i(u) \leq 0 \quad \forall i, h_j(u) = 0 \quad \forall j\}$$

$$f, g_1, \dots, g_r, h_1, \dots, h_s : \mathbb{R}^d \rightarrow \mathbb{R}$$

Lagrangian

- The *Lagrangian* is

$$L(u, \lambda, \nu) = f(u) + \sum_{i=1}^r \lambda_i g_i(u) + \sum_{j=1}^s \nu_j h_j(u)$$

- $\lambda = [\lambda_1, \dots, \lambda_r]^T$ and $\nu = [\nu_1, \dots, \nu_s]^T$ are called

Lagrange multipliers, dual variables

Dual Function

- The *Lagrangian dual function* is

$$L_D(\lambda, \nu) = \min_{u \in \mathbb{R}^d} L(u, \lambda, \nu)$$

- L_D is always concave
- The *dual optimization problem* is

$$\max_{\lambda, \nu: \lambda_i \geq 0} L_D(\lambda, \nu) = \max_{\lambda, \nu: \lambda_i \geq 0} \left[\min_{u \in \mathbb{R}^d} L(u, \lambda, \nu) \right]$$

- The original constrained optimization problem is sometimes called the

primal (optimization) problem

Rewriting the Primal

- Recall the Lagrangian:

$$L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu}) := f(\mathbf{u}) + \sum_{i=1}^r \lambda_i g_i(\mathbf{u}) + \sum_{j=1}^s \nu_j h_j(\mathbf{u})$$

- Observe that the primal may be re-written

$$\min_{\mathbf{u}} \left[\max_{\mathbf{u}, \boldsymbol{\lambda} : \lambda_i \geq 0} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu}) \right]$$
$$= \begin{cases} f(\mathbf{u}) \\ \infty \end{cases}$$

\mathbf{u} feasible

\mathbf{u} not feasible

Weak Duality

- Denote the optimal objective function values of the primal and dual

$$p^* = \min_{\mathbf{u}} \max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu})$$

$$d^* = \max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} \min_{\mathbf{u}} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = \max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} L_D(\boldsymbol{\lambda}, \boldsymbol{\nu})$$

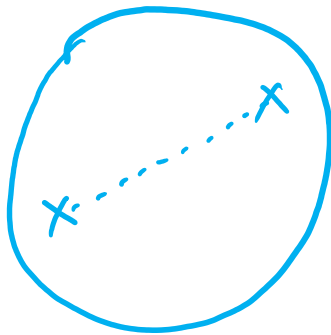
- *Weak duality* refers to the following fact, which always holds.
- **Theorem:** $d^* \leq p^*$

Poll

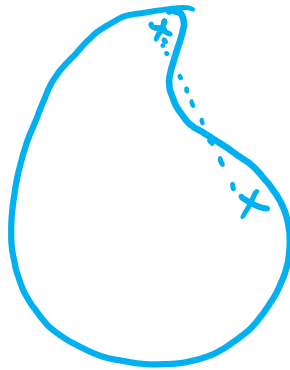
True or False: The constraint set (the set of all feasible points) is a *convex set* iff all of the functions $g_1, \dots, g_r, h_1, \dots, h_s$ are *convex functions*.

(A) True

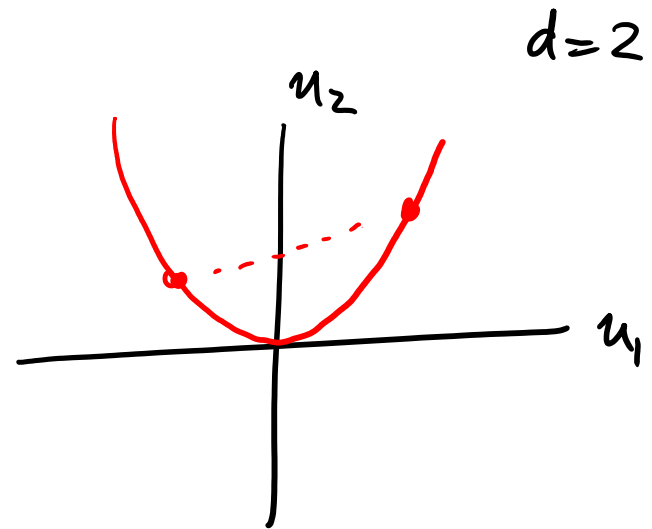
(B) False



convex set



not a
convex set



$$h_1(u) = -u_2 + u_1^2$$

Strong Duality

- If $p^* = d^*$, we say *strong duality* holds.
- The original constrained optimization problem is said to be *convex* if f and g_1, \dots, g_r are convex functions and h_1, \dots, h_s are affine.
- We state the following without proof.
- **Theorem:** If the original problem is convex and a constraint qualification holds, then $p^* = d^*$.
- Examples of constraint qualifications:
 - All g_i are affine
 - (Strict feasibility) $\exists \mathbf{u}$ s.t. $h_j(\mathbf{u}) = 0 \ \forall j$ and $\underline{g_i(\mathbf{u})} < 0 \ \forall i$

Affine: An affine function is a function of the form $\mathbf{u} \mapsto \mathbf{c}^T \mathbf{u} + d$

Poll

- Recall the constrained optimization problem defining the optimal soft-margin hyperplane classifier:

$$u = \begin{bmatrix} w \\ b \\ \xi \end{bmatrix}$$

$$\min_{w, b, \xi} \quad \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i = f(u) \quad (\text{OSM})$$

$$\text{s.t.} \quad \left. \begin{array}{l} y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i \\ \xi_i \geq 0 \quad \forall i \end{array} \right\} g_i(u), \quad i \in \{1, \dots, 2n\}$$

- True or false: Strong duality holds for problem OSM.
(A) True
(B) False

Big Picture

- For unconstrained optimization problems with differentiable objective function, we saw that

- $\nabla f(\mathbf{u}) = \mathbf{0}$ is *necessary* for \mathbf{u} to be a global minimizer
- If f is convex, then $\nabla f(\mathbf{u}) = \mathbf{0}$ is *sufficient* for \mathbf{u} to be a global minimizer

- For constrained optimization problems with differentiable objective and constraints, a similar result holds where $\nabla f(\mathbf{u}) = \mathbf{0}$ is replaced by the

Karush - Kuhn - Tucker (KKT) conditions

- We can use these conditions to solve and understand constrained optimization problems.

KKT Conditions: Necessity

- From now on assume f , g_i and h_j are all differentiable.
- Theorem:** If $p^* = d^*$, \mathbf{u}^* is primal optimal, and $(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$ is dual optimal, then the KKT conditions hold:

zero
vector

$$1) \nabla_{\mathbf{u}} f(\mathbf{u}^*) + \sum_{i=1}^r \lambda_i^* \nabla_{\mathbf{u}} g_i(\mathbf{u}^*) + \sum_{j=1}^s \nu_j^* \nabla_{\mathbf{u}} h_j(\mathbf{u}^*) = \mathbf{0} \quad \downarrow$$

$$2) g_i(\mathbf{u}^*) \leq 0 \quad \forall i$$

$$3) h_j(\mathbf{u}^*) = 0 \quad \forall j$$

$$4) \lambda_i^* \geq 0 \quad \forall i$$

$$5) \lambda_i^* g_i(\mathbf{u}^*) = 0 \quad \forall i$$

(complementary slackness)

KKT Conditions: Sufficiency

- **Theorem:** If the original problem is convex and $\tilde{\mathbf{u}}$, $\tilde{\boldsymbol{\lambda}}$, $\tilde{\boldsymbol{\nu}}$ satisfy the KKT conditions

1. $\nabla_{\mathbf{u}} f(\tilde{\mathbf{u}}) + \sum_{i=1}^r \tilde{\lambda}_i \nabla_{\mathbf{u}} g_i(\tilde{\mathbf{u}}) + \sum_{j=1}^s \tilde{\nu}_j \nabla_{\mathbf{u}} h_j(\tilde{\mathbf{u}}) = \mathbf{0}$
2. $g_i(\tilde{\mathbf{u}}) \leq 0 \ \forall i$
3. $h_j(\tilde{\mathbf{u}}) = 0 \ \forall j$
4. $\tilde{\lambda}_i \geq 0 \ \forall i$
5. $\tilde{\lambda}_i g_i(\tilde{\mathbf{u}}) = 0 \ \forall i$

then $\tilde{\mathbf{u}}$ is primal optimal, $(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$ is dual optimal, and strong duality holds.

How Is This Useful?

- Sometimes it is easier to solve the primal (analytically or computationally) by first solving the dual, and then using the KKT conditions to relate the primal solution to the dual solution.
- For example: suppose strong duality holds. If (λ^*, ν^*) is dual optimal, then (by the KKT necessity theorem) any primal optimal point u^* is a solution of

$$\nabla_u f(u^*) + \sum_{i=1}^r \lambda_i^* \nabla_{u_i} g_i(u^*) + \sum_{j=1}^s \nu_j^* \nabla_u h_j(u^*) = 0$$

- We can use this to find a u^* . In other words, we can recover a primal solution from a dual solution.

OSM Hyperplane

- Let's apply this theory to the OSM hyperplane classifier.

$$\min_{\mathbf{w}, b, \xi} \quad \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \quad (\text{OSM})$$

$$\text{s.t.} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i \in \{1, \dots, n\} \quad (\alpha_i)$$

$$\xi_i \geq 0 \quad \forall i \in \{1, \dots, n\}$$

$$\mathbf{u} = \begin{bmatrix} \mathbf{w} \\ b \\ \xi \end{bmatrix} \in \mathbb{R}^{d+1+n}, \quad \boldsymbol{\xi} = \begin{bmatrix} \xi_1 \\ \vdots \\ \xi_n \end{bmatrix}, \quad \boldsymbol{\lambda} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} = \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix}$$

- First observation: This is a convex optimization problem, and the constraint functions are all affine, hence strong duality holds.

Lagrangian

The Lagrangian is

$$\begin{aligned}
 L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) &= \overbrace{\frac{1}{2} \|\mathbf{w}\|^2}^{f(\mathbf{w})} + \frac{C}{n} \sum_{i=1}^n \xi_i - \underbrace{\sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i)}_{-g_i(\mathbf{w}), i=1, \dots, n} - \sum_{i=1}^n \beta_i \xi_i \\
 &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \beta_i \xi_i \\
 &\quad - b \sum_{i=1}^n \alpha_i y_i + \sum_{i=1}^n \xi_i \left(\frac{C}{n} - \alpha_i - \beta_i \right).
 \end{aligned}$$

$-g_i(\mathbf{w}), i=n+1, \dots, 2n$

Dual Function

- The dual function is

$$L_D(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \min_{\boldsymbol{w}, b, \boldsymbol{\xi}} L(\boldsymbol{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}).$$

- The optimization problem defining the dual function is an unconstrained minimization with a convex, differentiable objective function.
- Therefore, for fixed $\boldsymbol{\alpha}, \boldsymbol{\beta}$, we know that \boldsymbol{w}, b and $\boldsymbol{\xi}$ achieve the minimum iff

$$\left. \begin{aligned} \frac{\partial L}{\partial \boldsymbol{w}} &= \boldsymbol{w} - \sum \alpha_i y_i \boldsymbol{x}_i = \mathbf{0} \\ \frac{\partial L}{\partial b} &= - \sum \alpha_i y_i = 0 \\ \frac{\partial L}{\partial \xi_i} &= \frac{C}{n} - \alpha_i - \beta_i = 0 \quad \forall i. \end{aligned} \right\}$$

Dual Function

- Plugging in these formulas the Lagrangian simplifies to

$$\begin{aligned} L_D(\boldsymbol{\alpha}, \boldsymbol{\beta}) &= \frac{1}{2} \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_i \alpha_i y_i \left\langle \sum_j \alpha_j y_j \mathbf{x}_j, \mathbf{x}_i \right\rangle + \sum_i \alpha_i \\ &= \frac{1}{2} \left\langle \sum_i \alpha_i y_i \mathbf{x}_i, \sum_j \alpha_j y_j \mathbf{x}_j \right\rangle - \left\langle \sum_i \alpha_i y_i \mathbf{x}_i, \sum_j \alpha_j y_j \mathbf{x}_j \right\rangle + \sum_i \alpha_i \\ &= -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \alpha_i. \end{aligned}$$

Dual Optimization Problem

- Therefore, the dual optimization problem

$$\max_{\alpha \geq 0, \beta \geq 0} L_D(\alpha, \beta)$$

may be written

$$\max_{\alpha, \beta} -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^n \alpha_i$$

$$\text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i = 0$$

$$\alpha_i + \beta_i = \frac{C}{n} \quad \forall i$$

$$\alpha_i \geq 0, \beta_i \geq 0 \quad \forall i.$$

$$\Leftrightarrow 0 \leq \alpha_i \leq \frac{C}{n} \quad \forall i$$

Dual Optimization Problem

- In summary, the dual optimization problem can be written

$$\begin{aligned} \max_{\alpha} \quad & -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^n \alpha_i \\ \text{s.t} \quad & \sum_i \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq \frac{C}{n}, \quad \forall i = 1, \dots, n \end{aligned}$$

- This is another quadratic program.
- What do you notice? *dot products*

Solving Primal From Dual

after

- The problem is convex with ~~linear~~ constraints so strong duality holds.
- Let α^* denote a dual solution, obtained by solving the dual QP.
- Let (w^*, b^*, ξ^*) denote a primal solution; we don't know it yet, just suppose it exists
- By the KKT necessity theorem, $(w^*, b^*, \xi^*, \alpha^*)$ satisfy the KKT conditions.
- We may use these to compute (w^*, b^*, ξ^*) from α^* .
- From the first KKT condition

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

- See lecture notes for how to recover b^* .

Final Classifier

$$\begin{aligned}h(x) &= \text{sign} \{ (w^*)^T x + b^* \} \\ &= \text{sign} \left\{ \sum_{i=1}^n \alpha_i^* y_i \langle x_i, x \rangle + b^* \right\}\end{aligned}$$

Summary: Dual Formulation

- The OSM hyperplane is kernelizable: The dual is

$$\begin{aligned} \max_{\alpha} \quad & -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^n \alpha_i \\ \text{s.t} \quad & \sum_i \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq \frac{C}{n}, \quad \forall i = 1, \dots, n \end{aligned}$$

- The classifier is expressed

$$f(\mathbf{x}) = \text{sign} \left\{ \sum_{i=1}^n \alpha_i^* y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b^* \right\}$$

- Taking any j with $0 < \alpha_j^* < \frac{C}{n}$, the offset is given by

$$b^* = y_j - \sum_{i=1}^n \alpha_i^* y_i \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

The Support Vector Machine

- Let k be an inner product/SPD kernel
- The *support vector machine* is the classifier obtained by solving

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i \\ \text{s.t} \quad & \sum_i \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq \frac{C}{n}, \quad \forall i = 1, \dots, n \end{aligned}$$

- The classifier is expressed

$$f(\mathbf{x}) = \text{sign} \left\{ \sum_{i=1}^n \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}) + b^* \right\}$$

where $b^* = y_j - \sum_{i=1}^n \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}_j)$ for any j with $0 < \alpha_j^* < \frac{C}{n}$.

Support Vectors

- From complementary slackness,

- If \mathbf{x}_i satisfies

we call \mathbf{x}_i a *support vector*.

- Therefore, if \mathbf{x}_i is *not* a support vector, then
- Conclusion: \mathbf{w}^* depends only on the SVs: