

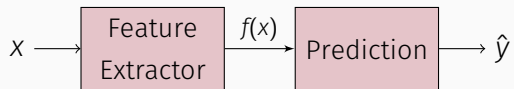
KERNEL SUBSPACE AND FEATURE EXTRACTION

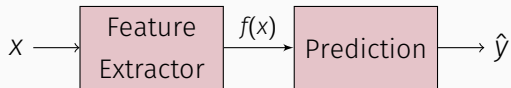
Xiangxiang Xu

Joint work with Prof. Lihong Zheng

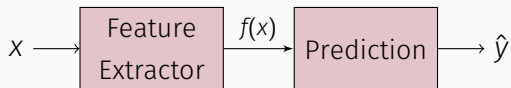
ISIT 2023





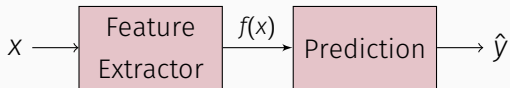


Features: Explicit & Implicit



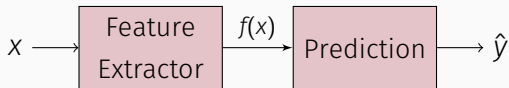
Features: Explicit & Implicit

- ▶ Parameterized features, e.g., deep neural networks



Features: Explicit & Implicit

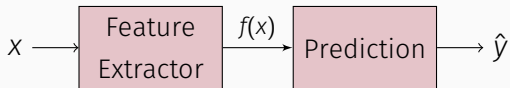
- ▶ Parameterized features, e.g., deep neural networks
- ▶ Kernel methods: features defined by kernel



Features: Explicit & Implicit

- ▶ Parameterized features, e.g., deep neural networks
- ▶ Kernel methods: features defined by kernel

Understand Kernel

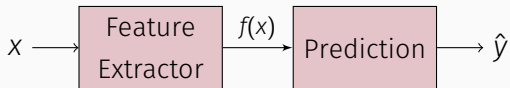


Features: Explicit & Implicit

- ▶ Parameterized features, e.g., deep neural networks
- ▶ Kernel methods: features defined by kernel

Understand Kernel

- ▶ How it worked?

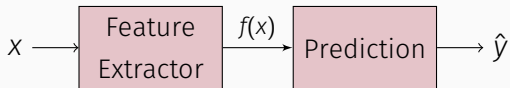


Features: Explicit & Implicit

- ▶ Parameterized features, e.g., deep neural networks
- ▶ Kernel methods: features defined by kernel

Understand Kernel

- ▶ How it worked?
- ▶ What makes a kernel good?



Features: Explicit & Implicit

- ▶ Parameterized features, e.g., deep neural networks
- ▶ Kernel methods: features defined by kernel

Understand Kernel

- ▶ How it worked?
- ▶ What makes a kernel good?
- ▶ How to obtain a good kernel?

Data X , Label $Y = \pm 1$, balanced

Data X , Label $Y = \pm 1$, balanced

$$\text{MAP} \equiv \text{ML}: \quad P_{X|Y=1}(x) \underset{\hat{y}=-1}{\overset{\hat{y}=1}{\geq}} P_{X|Y=-1}(x)$$

Data X , Label $Y = \pm 1$, balanced

$$\text{MAP} \equiv \text{ML}: \quad P_{X|Y=1}(x) \underset{\hat{y}=-1}{\overset{\hat{y}=1}{\geq}} P_{X|Y=-1}(x)$$

A Feature Perspective

$$\text{Consider } f^*(x) \propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)}$$

- ▶ $\mathbb{E}[f^*(X)] = 0$, $\text{var}(f^*(X)) = 1$ (pick the scaling factor)

Data X , Label $Y = \pm 1$, balanced

$$\text{MAP} \equiv \text{ML}: \quad P_{X|Y=1}(x) \underset{\hat{y}=-1}{\overset{\hat{y}=1}{\geq}} P_{X|Y=-1}(x)$$

A Feature Perspective

$$\text{Consider } f^*(x) \propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)}$$

- ▶ $\mathbb{E}[f^*(X)] = 0$, $\text{var}(f^*(X)) = 1$ (pick the scaling factor)
- ▶ maximally correlated choice

Data X , Label $Y = \pm 1$, balanced

$$\text{MAP} \equiv \text{ML}: \quad P_{X|Y=1}(x) \underset{\hat{y}=-1}{\overset{\hat{y}=1}{\geq}} P_{X|Y=-1}(x)$$

A Feature Perspective

$$\text{Consider } f^*(x) \propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)}$$

- ▶ $\mathbb{E}[f^*(X)] = 0$, $\text{var}(f^*(X)) = 1$ (pick the scaling factor)
- ▶ maximally correlated choice
- ▶ $\frac{P_{Y|X}(y|x)}{P_Y(y)} = 1 + \varrho \cdot y \cdot f^*(x)$ ϱ : maximal correlation

Data X , Label $Y = \pm 1$, balanced

$$\text{MAP} \equiv \text{ML}: \quad P_{X|Y=1}(x) \underset{\hat{y}=-1}{\overset{\hat{y}=1}{\geq}} P_{X|Y=-1}(x)$$

A Feature Perspective

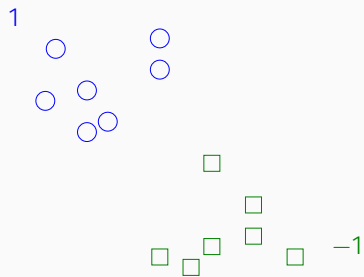
$$\text{Consider } f^*(x) \propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)}$$

- ▶ $\mathbb{E}[f^*(X)] = 0$, $\text{var}(f^*(X)) = 1$ (pick the scaling factor)
- ▶ maximally correlated choice
- ▶ $\frac{P_{Y|X}(y|x)}{P_Y(y)} = 1 + \varrho \cdot y \cdot f^*(x)$ ϱ : maximal correlation

Optimal Decision

$$\hat{y}_{\text{MAP}}(x) = \text{sgn}(f^*(x)) \quad f^*: \text{ sufficient statistic}$$

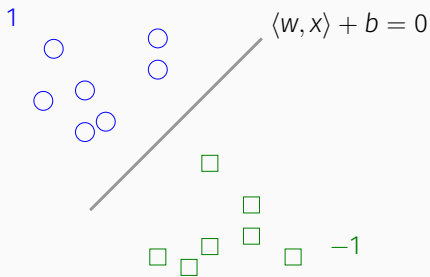
SUPPORT VECTOR MACHINE



SUPPORT VECTOR MACHINE

Separating Hyperplane

$$\langle w, x \rangle + b = 0$$

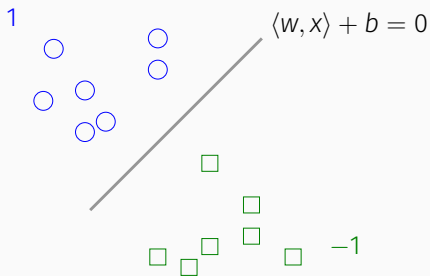


SUPPORT VECTOR MACHINE

Separating Hyperplane

$$\langle w, x \rangle + b = 0$$

Decide which side x lies in:



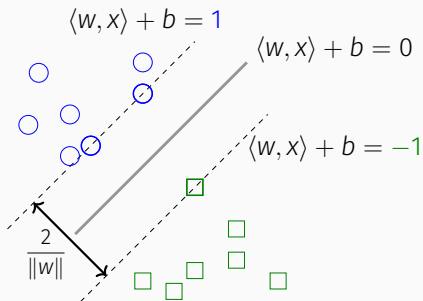
Separating Hyperplane

$$\langle w, x \rangle + b = 0$$

Decide which side x lies in:

$$\hat{y}_{\text{SVM}}(x) \triangleq \text{sgn}(\langle w^*, x \rangle + b^*)$$

w^*, b^* : maximize margin



Separating Hyperplane

$$\langle w, x \rangle + b = 0$$

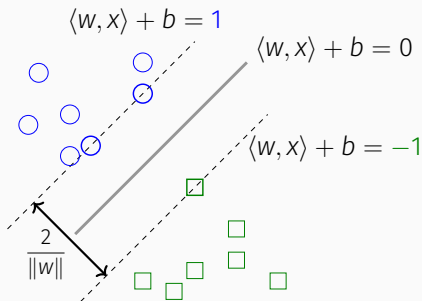
Decide which side x lies in:

$$\hat{y}_{\text{SVM}}(x) \triangleq \text{sgn}(\langle w^*, x \rangle + b^*)$$

w^*, b^* : maximize margin

$$\min_{w, b} \frac{1}{2} \cdot \|w\|^2$$

$$\text{s. t. } 1 - y_i \cdot (\langle w, x_i \rangle + b) \leq 0$$



Separating Hyperplane

$$\langle w, x \rangle + b = 0$$

Decide which side x lies in:

$$\hat{y}_{\text{SVM}}(x) \triangleq \text{sgn}(\langle w^*, x \rangle + b^*)$$

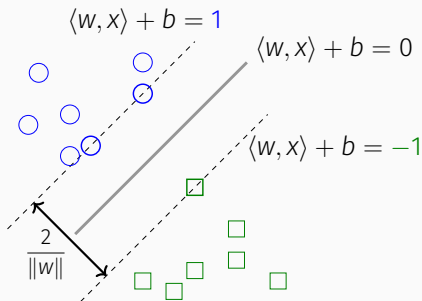
w^*, b^* : maximize margin

$$\min_{w, b} \frac{1}{2} \cdot \|w\|^2$$

$$\text{s. t. } 1 - y_i \cdot (\langle w, x_i \rangle + b) \leq 0$$

Soft Margin

$$L_{\text{SVM}}(w, b; \lambda) \triangleq \mathbb{E} [(1 - Y \cdot (\langle w, X \rangle + b))^+] + \frac{\lambda}{2} \cdot \|w\|^2$$



Implementation

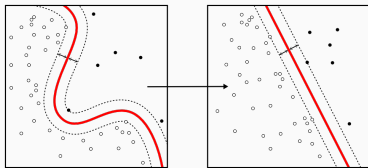
- ▶ solve the dual optimization problem
- ▶ expressed in terms of $\langle x_i, x_j \rangle$

Implementation

- ▶ solve the dual optimization problem
- ▶ expressed in terms of $\langle x_i, x_j \rangle$

Introducing Features

- ▶ use $f(X)$ instead of X
- ▶ linear separability
- ▶ “flexibility”

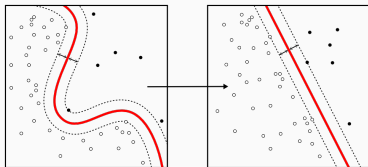


Implementation

- ▶ solve the dual optimization problem
- ▶ expressed in terms of $\langle x_i, x_j \rangle$

Introducing Features

- ▶ use $f(X)$ instead of X
- ▶ linear separability
- ▶ “flexibility”



Kernel $k(x, x') = \langle f(x), f(x') \rangle$

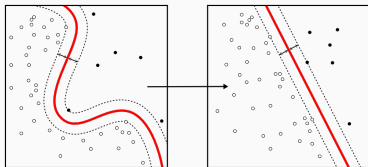
- ▶ resulting f : implicit and high-dimensional
- ▶ map data without extra computational cost

Implementation

- ▶ solve the dual optimization problem
- ▶ expressed in terms of $\langle x_i, x_j \rangle$

Introducing Features

- ▶ use $f(X)$ instead of X
- ▶ linear separability
- ▶ “flexibility”



Kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

- ▶ resulting f : implicit and high-dimensional
- ▶ map data without extra computational cost
- ▶ choices of κ : RBF, polynomial

“Soft” Regime: large λ

$$w^* = \lambda^{-1} \cdot \mathbb{E} [f(X)Y], \quad b^* = 0$$

– assume $\mathbb{E} [f] = 0, \text{cov}(f) = I$

“Soft” Regime: large λ

$$w^* = \lambda^{-1} \cdot \mathbb{E}[f(X)Y], \quad b^* = 0 \quad - \text{assume } \mathbb{E}[f] = 0, \text{cov}(f) = I$$

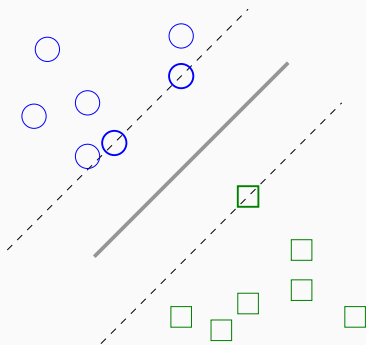
$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle w^*, f(x) \rangle + b^*) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

“Soft” Regime: large λ

$$w^* = \lambda^{-1} \cdot \mathbb{E}[f(X)Y], \quad b^* = 0 \quad - \text{assume } \mathbb{E}[f] = 0, \text{cov}(f) = I$$

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle w^*, f(x) \rangle + b^*) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

$$\mathbb{E}[f(X)Y] \propto \mathbb{E}[f(X)|Y=1] - \mathbb{E}[f(X)|Y=-1]$$

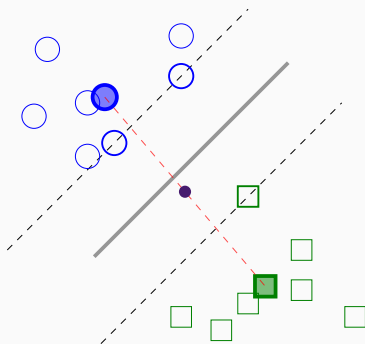


“Soft” Regime: large λ

$$w^* = \lambda^{-1} \cdot \mathbb{E}[f(X)Y], \quad b^* = 0 \quad - \text{assume } \mathbb{E}[f] = 0, \text{cov}(f) = I$$

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle w^*, f(x) \rangle + b^*) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

$$\mathbb{E}[f(X)Y] \propto \mathbb{E}[f(X)|Y=1] - \mathbb{E}[f(X)|Y=-1]$$

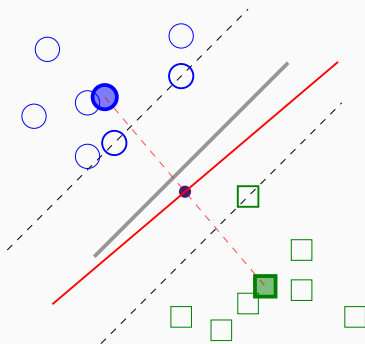


“Soft” Regime: large λ

$$w^* = \lambda^{-1} \cdot \mathbb{E}[f(X)Y], \quad b^* = 0 \quad - \text{assume } \mathbb{E}[f] = 0, \text{cov}(f) = I$$

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle w^*, f(x) \rangle + b^*) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

$$\mathbb{E}[f(X)Y] \propto \mathbb{E}[f(X)|Y=1] - \mathbb{E}[f(X)|Y=-1]$$



SVM FOR GIVEN f

“Soft” Regime: large λ

$$w^* = \lambda^{-1} \cdot \mathbb{E}[f(X)Y], \quad b^* = 0 \quad - \text{assume } \mathbb{E}[f] = 0, \text{cov}(f) = I$$

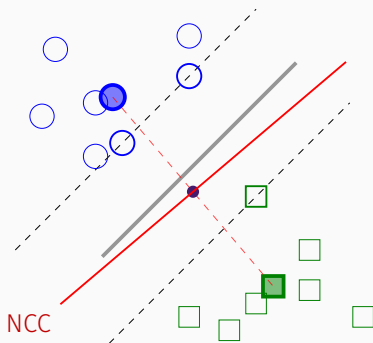
$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle w^*, f(x) \rangle + b^*) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

$$\mathbb{E}[f(X)Y] \propto \mathbb{E}[f(X)|Y = 1] - \mathbb{E}[f(X)|Y = -1]$$

Equivalent Form

$$\arg \min_y \|f(x) - \mathbb{E}[f(X)|Y = y]\|$$

- ▶ Nearest Centroid Classifier



SVM FOR GIVEN f

“Soft” Regime: large λ

$$w^* = \lambda^{-1} \cdot \mathbb{E}[f(X)Y], \quad b^* = 0 \quad - \text{assume } \mathbb{E}[f] = 0, \text{cov}(f) = I$$

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle w^*, f(x) \rangle + b^*) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

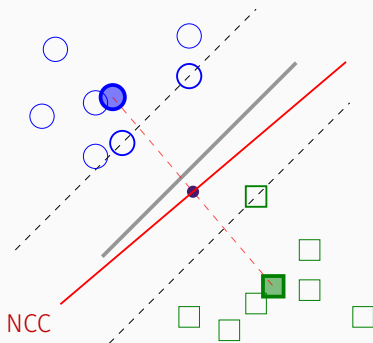
$$\mathbb{E}[f(X)Y] \propto \mathbb{E}[f(X)|Y=1] - \mathbb{E}[f(X)|Y=-1]$$

Equivalent Form

$$\arg \min_y \|\mathbb{E}[f(X)|Y=y]\|$$

- Nearest Centroid Classifier

How to pick f ?



SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

- ▶ Based on a linear combination of f_1, \dots, f_d

SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

- ▶ Based on a linear combination of f_1, \dots, f_d

Optimal Decision

$$\hat{y}_{\text{MAP}}(x) = \text{sgn}(f^*(x)) \quad f^*: \text{max. corr. function}$$

SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

- ▶ Based on a linear combination of f_1, \dots, f_d

Optimal Decision

$$\hat{y}_{\text{MAP}}(x) = \text{sgn}(f^*(x)) \quad f^*: \text{max. corr. function}$$

- ▶ SVM is optimal if $f = f^*$ ($d = 1$)

SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

- ▶ Based on a linear combination of f_1, \dots, f_d

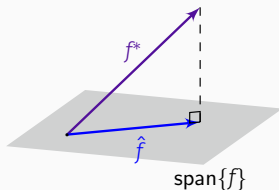
Optimal Decision

$$\hat{y}_{\text{MAP}}(x) = \text{sgn}(f^*(x)) \quad f^*: \text{max. corr. function}$$

- ▶ SVM is optimal if $f = f^*$ ($d = 1$)

Optimal Linear Combination

\hat{f} : project f^* onto $\text{span}\{f_1, \dots, f_d\}$



SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$$

- ▶ Based on a linear combination of f_1, \dots, f_d

Optimal Decision

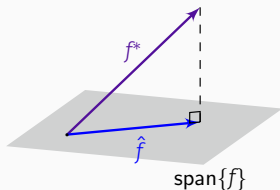
$$\hat{y}_{\text{MAP}}(x) = \text{sgn}(f^*(x)) \quad f^*: \text{max. corr. function}$$

- ▶ SVM is optimal if $f = f^*$ ($d = 1$)

Optimal Linear Combination

\hat{f} : project f^* onto $\text{span}\{f_1, \dots, f_d\}$

- ▶ Feature space $\mathcal{F} \triangleq \{\phi: \mathcal{X} \rightarrow \mathbb{R}\}$
- ▶ Inner product $\langle \phi_1, \phi_2 \rangle_{\mathcal{F}} \triangleq \mathbb{E}[\phi_1(X)\phi_2(X)]$ (data dependent)



SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle) = \text{sgn}(\hat{f}(x))$$

- ▶ Based on a linear combination of f_1, \dots, f_d

Optimal Decision

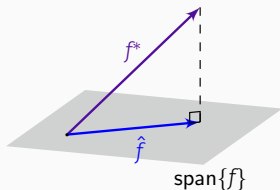
$$\hat{y}_{\text{MAP}}(x) = \text{sgn}(f^*(x)) \quad f^*: \text{max. corr. function}$$

- ▶ SVM is optimal if $f = f^*$ ($d = 1$)

Optimal Linear Combination

\hat{f} : project f^* onto $\text{span}\{f_1, \dots, f_d\}$

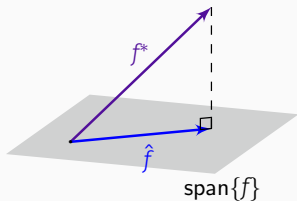
- ▶ Feature space $\mathcal{F} \triangleq \{\phi: \mathcal{X} \rightarrow \mathbb{R}\}$
- ▶ Inner product $\langle \phi_1, \phi_2 \rangle_{\mathcal{F}} \triangleq \mathbb{E}[\phi_1(X)\phi_2(X)]$ (data dependent)



SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\hat{f}(x))$$

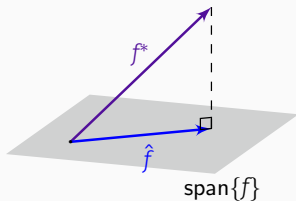
► $\hat{f} = \Pi(f^*; \text{span}\{f\})$



SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\hat{f}(x))$$

► $\hat{f} = \Pi(f^*; \text{span}\{f\})$



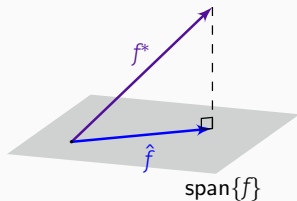
Optimal Decision

$$\hat{f} = f^* \implies \hat{y}_{\text{SVM}} = \hat{y}_{\text{MAP}}$$

SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\hat{f}(x))$$

► $\hat{f} = \Pi(f^*; \text{span}\{f\})$



Optimal Decision

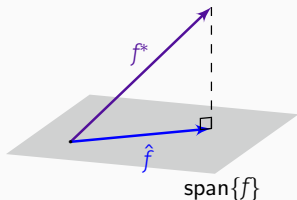
$$\hat{f} = f^* \implies \hat{y}_{\text{SVM}} = \hat{y}_{\text{MAP}}$$

► require $f^* \in \text{span}\{f\}$

SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\hat{f}(x))$$

► $\hat{f} = \Pi(f^*; \text{span}\{f\})$



Optimal Decision

$$\hat{f} = f^* \implies \hat{y}_{\text{SVM}} = \hat{y}_{\text{MAP}}$$

► require $f^* \in \text{span}\{f\}$

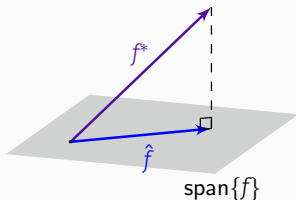
SVM Loss

$$\text{Put } w^*, b^* \text{ in } \implies L_{\text{SVM}}^*(f) = 1 - \frac{c}{\lambda} \cdot \|\hat{f}\|_{\mathcal{F}}^2$$

SVM Decision

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\hat{f}(x))$$

► $\hat{f} = \Pi(f^*; \text{span}\{f\})$



Optimal Decision

$$\hat{f} = f^* \implies \hat{y}_{\text{SVM}} = \hat{y}_{\text{MAP}}$$

► require $f^* \in \text{span}\{f\}$

SVM Loss

Put w^*, b^* in $\implies L_{\text{SVM}}^*(f) = 1 - \frac{c}{\lambda} \cdot \|\hat{f}\|_{\mathcal{F}}^2$

► “oracle” feature extractor: minimize $L_{\text{SVM}}^*(f)$ over f

Decision in terms of kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

Decision in terms of kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

General Form

Rewrite $\hat{y}_{\text{SVM}}(x; f) = \text{sgn} (\langle \mathbb{E}[f(X)Y], f(x) \rangle)$

Decision in terms of kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

General Form

Rewrite $\hat{y}_{\text{SVM}}(x; f) = \text{sgn} (\langle \mathbb{E}[f(X)Y], f(x) \rangle)$

Decision in terms of kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

General Form

Rewrite $\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) \cdot Y])$$

Decision in terms of kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

General Form

Rewrite $\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) \cdot Y])$$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) f^*(X)])$$

— Replace P_{XY} by $P_X P_Y$ and f^*

Decision in terms of kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

General Form

Rewrite $\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) \cdot Y])$$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) f^*(X)])$$

— Replace P_{XY} by $P_X P_Y$ and f^*

Key Quantity

$$\mathbb{E}[\kappa(X, x) f^*(X)]$$

Decision in terms of kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

General Form

Rewrite $\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) \cdot Y])$$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) f^*(X)]) \quad - \text{Replace } P_{XY} \text{ by } P_X P_Y \text{ and } f^*$$

Key Quantity

$$\mathbb{E}[\kappa(X, x) f^*(X)] = \int \kappa(x, x') \cdot f^*(x') p_X(x') dx'$$

Decision in terms of kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

General Form

Rewrite $\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) \cdot Y])$$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) f^*(X)]) \quad - \text{Replace } P_{XY} \text{ by } P_X P_Y \text{ and } f^*$$

Key Quantity

$$\mathbb{E}[\kappa(X, x) f^*(X)] = \int \kappa(x, x') \cdot f^*(x') p_X(x') dx' = [\kappa \circ f^*](x)$$

– “kernel” of integral operator

Decision in terms of kernel $\kappa(x, x') = \langle f(x), f(x') \rangle$

General Form

Rewrite $\hat{y}_{\text{SVM}}(x; f) = \text{sgn}(\langle \mathbb{E}[f(X)Y], f(x) \rangle)$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) \cdot Y])$$

$$= \text{sgn}(\mathbb{E}[\kappa(X, x) f^*(X)]) \quad - \text{Replace } P_{XY} \text{ by } P_X P_Y \text{ and } f^*$$

Key Quantity

$$\mathbb{E}[\kappa(X, x) f^*(X)] = \int \kappa(x, x') \cdot f^*(x') p_X(x') dx' = [\kappa \circ f^*](x)$$

– “kernel” of integral operator

How to pick κ ?

Kernel for Feature Subspace \mathcal{G}

$$\kappa_{\mathcal{G}}(x, x') \triangleq \langle \phi(x), \phi(x') \rangle \quad \phi = (\phi_1, \dots, \phi_d)^T: \text{ orthonormal basis}$$

Kernel for Feature Subspace \mathcal{G}

$\kappa_{\mathcal{G}}(x, x') \triangleq \langle \phi(x), \phi(x') \rangle$ $\phi = (\phi_1, \dots, \phi_d)^T$: orthonormal basis

► $\kappa_{\mathcal{G}} \circ f = \Pi(f; \mathcal{G})$

— Projection onto \mathcal{G}

Kernel for Feature Subspace \mathcal{G}

$\kappa_{\mathcal{G}}(x, x') \triangleq \langle \phi(x), \phi(x') \rangle$ $\phi = (\phi_1, \dots, \phi_d)^T$: orthonormal basis

▶ $\kappa_{\mathcal{G}} \circ f = \Pi(f; \mathcal{G})$

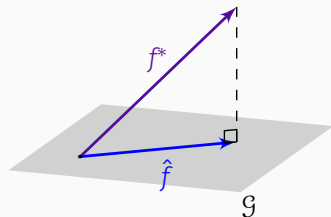
– Projection onto \mathcal{G}

SVM With $\kappa_{\mathcal{G}}$

▶ Decision $\hat{y}_{\text{SVM}}(x; \kappa_{\mathcal{G}}) = \text{sgn}(\hat{f}(x))$

▶ Loss $L_{\text{SVM}}^*(\kappa_{\mathcal{G}}) = 1 - \frac{c}{\lambda} \cdot \|\hat{f}\|_{\mathcal{F}}^2$

▶ $\hat{f} \triangleq \Pi(f^*; \mathcal{G})$



Kernel for Feature Subspace \mathcal{G}

$\kappa_{\mathcal{G}}(x, x') \triangleq \langle \phi(x), \phi(x') \rangle$ $\phi = (\phi_1, \dots, \phi_d)^T$: orthonormal basis

▶ $\kappa_{\mathcal{G}} \circ f = \Pi(f; \mathcal{G})$

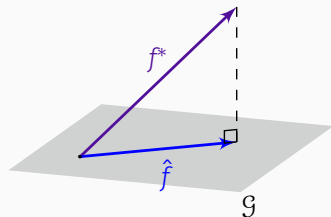
– Projection onto \mathcal{G}

SVM With $\kappa_{\mathcal{G}}$

▶ Decision $\hat{y}_{\text{SVM}}(x; \kappa_{\mathcal{G}}) = \text{sgn}(\hat{f}(x))$

▶ Loss $L_{\text{SVM}}^*(\kappa_{\mathcal{G}}) = 1 - \frac{c}{\lambda} \cdot \|\hat{f}\|_{\mathcal{F}}^2$

▶ $\hat{f} \triangleq \Pi(f^*; \mathcal{G})$



Quality of $\kappa_{\mathcal{G}}$

▶ Determined by $\|\hat{f}\|_{\mathcal{F}}$, or the angle between f^* and \mathcal{G}

Kernel for Feature Subspace \mathcal{G}

$\kappa_{\mathcal{G}}(x, x') \triangleq \langle \phi(x), \phi(x') \rangle$ $\phi = (\phi_1, \dots, \phi_d)^T$: orthonormal basis

▶ $\kappa_{\mathcal{G}} \circ f = \Pi(f; \mathcal{G})$

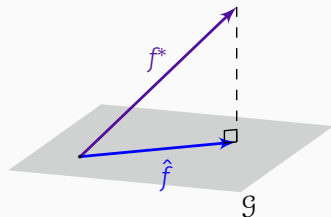
– Projection onto \mathcal{G}

SVM With $\kappa_{\mathcal{G}}$

▶ Decision $\hat{y}_{\text{SVM}}(x; \kappa_{\mathcal{G}}) = \text{sgn}(\hat{f}(x))$

▶ Loss $L_{\text{SVM}}^*(\kappa_{\mathcal{G}}) = 1 - \frac{c}{\lambda} \cdot \|\hat{f}\|_{\mathcal{F}}^2$

▶ $\hat{f} \triangleq \Pi(f^*; \mathcal{G})$



Quality of $\kappa_{\mathcal{G}}$

▶ Determined by $\|\hat{f}\|_{\mathcal{F}}$, or the angle between f^* and \mathcal{G}

▶ Optimal choice: $\mathcal{G} \ni f^*$

Two classes from the same family $\pi(x; \theta)$

▶ $Y = 1 \leftrightarrow \theta_0$

$$P_{X|Y=1}(x) = \pi(x; \theta_0)$$

▶ $Y = -1 \leftrightarrow -\theta_0$

$$P_{X|Y=-1}(x) = \pi(x; -\theta_0)$$

Two classes from the same family $\pi(x; \theta)$

▶ $Y = 1 \leftrightarrow \theta_0$

$$P_{X|Y=1}(x) = \pi(x; \theta_0)$$

▶ $Y = -1 \leftrightarrow -\theta_0$

$$P_{X|Y=-1}(x) = \pi(x; -\theta_0)$$

Optimal Kernel

$$f^*(x) \propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)}$$

Two classes from the same family $\pi(x; \theta)$

▶ $Y = 1 \leftrightarrow \theta_0$

$$P_{X|Y=1}(x) = \pi(x; \theta_0)$$

▶ $Y = -1 \leftrightarrow -\theta_0$

$$P_{X|Y=-1}(x) = \pi(x; -\theta_0)$$

Optimal Kernel

$$f^*(x) \propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)} \propto \frac{\nabla_{\theta} \pi(x; \theta)|_{\theta=0}}{\pi(x; 0)}$$

Two classes from the same family $\pi(x; \theta)$

▶ $Y = 1 \leftrightarrow \theta_0$

$$P_{X|Y=1}(x) = \pi(x; \theta_0)$$

▶ $Y = -1 \leftrightarrow -\theta_0$

$$P_{X|Y=-1}(x) = \pi(x; -\theta_0)$$

Optimal Kernel

$$\begin{aligned} f^*(x) &\propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)} \propto \frac{\nabla_{\theta} \pi(x; \theta)|_{\theta=0}}{\pi(x; 0)} \\ &= \nabla_{\theta} \log \pi(x; \theta)|_{\theta=0} \end{aligned}$$

Two classes from the same family $\pi(x; \theta)$

▶ $Y = 1 \leftrightarrow \theta_0$

$$P_{X|Y=1}(x) = \pi(x; \theta_0)$$

▶ $Y = -1 \leftrightarrow -\theta_0$

$$P_{X|Y=-1}(x) = \pi(x; -\theta_0)$$

Optimal Kernel

$$\begin{aligned} f^*(x) &\propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)} \propto \frac{\nabla_{\theta} \pi(x; \theta)|_{\theta=0}}{\pi(x; 0)} \\ &= \nabla_{\theta} \log \pi(x; \theta)|_{\theta=0} \end{aligned}$$

▶ f^* is score function $s(x)$

Two classes from the same family $\pi(x; \theta)$

▶ $Y = 1 \leftrightarrow \theta_0$

$$P_{X|Y=1}(x) = \pi(x; \theta_0)$$

▶ $Y = -1 \leftrightarrow -\theta_0$

$$P_{X|Y=-1}(x) = \pi(x; -\theta_0)$$

Optimal Kernel

$$\begin{aligned} f^*(x) &\propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)} \propto \frac{\nabla_{\theta} \pi(x; \theta)|_{\theta=0}}{\pi(x; 0)} \\ &= \nabla_{\theta} \log \pi(x; \theta)|_{\theta=0} \end{aligned}$$

▶ f^* is score function $s(x)$

▶ optimal $\kappa(x, x') = \langle s(x), s(x') \rangle$

Two classes from the same family $\pi(x; \theta)$

▶ $Y = 1 \leftrightarrow \theta_0$

$$P_{X|Y=1}(x) = \pi(x; \theta_0)$$

▶ $Y = -1 \leftrightarrow -\theta_0$

$$P_{X|Y=-1}(x) = \pi(x; -\theta_0)$$

Optimal Kernel

$$\begin{aligned} f^*(x) &\propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)} \propto \frac{\nabla_{\theta} \pi(x; \theta)|_{\theta=0}}{\pi(x; 0)} \\ &= \nabla_{\theta} \log \pi(x; \theta)|_{\theta=0} \end{aligned}$$

▶ f^* is score function $s(x)$

▶ optimal $\kappa(x, x') = \langle s(x), s(x') \rangle$

— Fisher kernel!

Two classes from the same family $\pi(x; \theta)$

▶ $Y = 1 \leftrightarrow \theta_0$

$$P_{X|Y=1}(x) = \pi(x; \theta_0)$$

▶ $Y = -1 \leftrightarrow -\theta_0$

$$P_{X|Y=-1}(x) = \pi(x; -\theta_0)$$

Optimal Kernel

$$\begin{aligned} f^*(x) &\propto \frac{P_{X|Y=1}(x) - P_{X|Y=-1}(x)}{P_X(x)} \propto \frac{\nabla_{\theta} \pi(x; \theta)|_{\theta=0}}{\pi(x; 0)} \\ &= \nabla_{\theta} \log \pi(x; \theta)|_{\theta=0} \end{aligned}$$

▶ f^* is score function $s(x)$

▶ optimal $\kappa(x, x') = \langle s(x), s(x') \rangle$

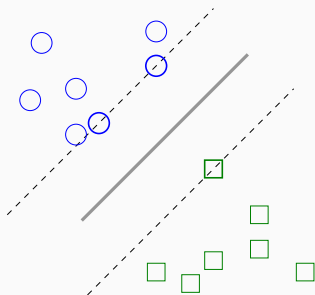
— Fisher kernel!

Quality of Kernel $\kappa_{\mathcal{G}}$

▶ $\|\Pi(s; \mathcal{G})\|_{\mathcal{F}}$

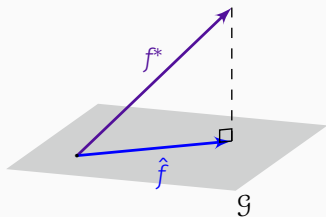
KERNEL SUBSPACE AND FEATURE EXTRACTION

Euclidean Geometry



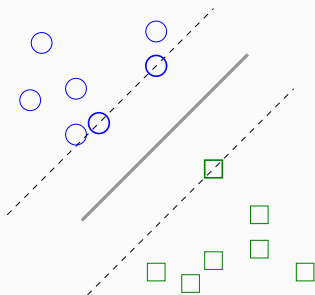
→

Feature Geometry



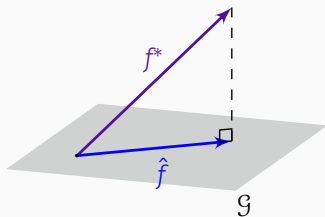
KERNEL SUBSPACE AND FEATURE EXTRACTION

Euclidean Geometry



→

Feature Geometry



- ▶ Feature subspace $\mathcal{G} \leftrightarrow$ Kernel $\kappa_{\mathcal{G}}$
- ▶ Performance decided by $\angle(f^*, \mathcal{G})$
- ▶ Quantitative metric of kernel quality