

S V M



Support Vector Machines

- Parametric classifier
- Motivation Geometric (maximizing the margin) Not probabilistic

Outline of our study

- 1) Two-class SVM
 - linearly separable
 - linearly nonseparable
 - "Kernel Trick"
- 2) Multi-class SVMs
 - aggregate / combined algorithms
 - One vs All (OvA)
 - One vs one (OvO)
 - uncombined algorithms (not covered)

Two-class Linearly Separable SVM

- Label space $\gamma = \{-1, +1\}$
- Feature vectors $X = \mathbb{R}^d$
- Training set $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$
- Family of classifiers: "linear" classifiers

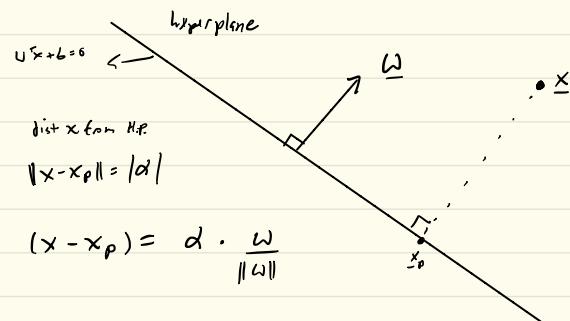
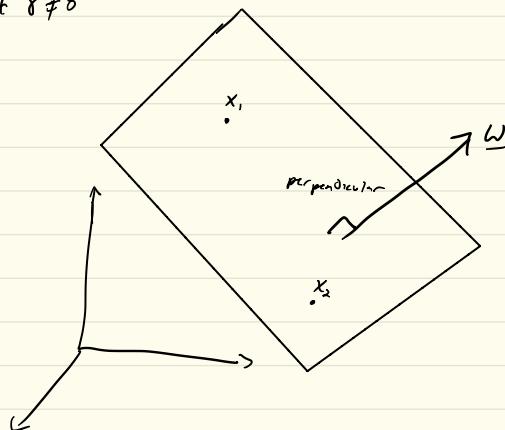
$$H = \{h(x) = \text{sign}(w^\top x + b) \quad w \in \mathbb{R}^d, \quad b \in \mathbb{R}\}$$

- Parameters
 $\Theta = \begin{pmatrix} w \\ b \end{pmatrix} \quad (d+1) \text{ scalar parameters}$

Equation of hyperplane in \mathbb{R}^d : $\omega^\top x + b = 0$

(ω, b) and $(\gamma\omega, \gamma b)$ define the same hyperplane if $\gamma \neq 0$

$$\begin{aligned} \omega^\top x_1 + b &= 0 \\ \omega^\top x_2 + b &= 0 \\ \omega^\top(x_1 - x_2) &= 0 \end{aligned}$$



\perp dist. (euclid.) of pt. x from hyperplane
 $= \frac{|\omega^\top x + b|}{\|\omega\|}$ ← euclidean

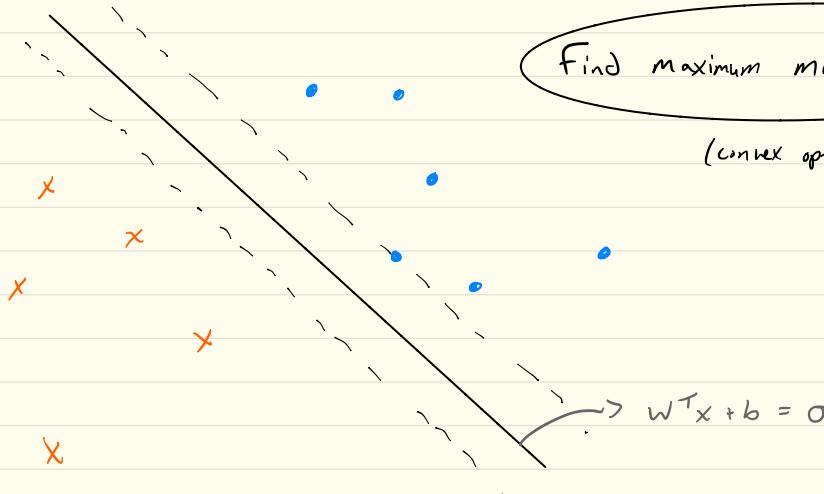
$$\omega^\top x_p + b = 0$$

$$\Rightarrow \omega^\top(x - x_p) = \alpha \frac{\omega^\top \omega}{\|\omega\|} = \alpha \frac{\|\omega\|^2}{\|\omega\|} = \alpha \|\omega\|$$

$$\Rightarrow \alpha = \frac{\omega^\top \omega - \omega^\top x_p}{\|\omega\|} = \frac{\omega^\top x + b}{\|\omega\|} \Rightarrow |\alpha| = \frac{|\omega^\top x + b|}{\|\omega\|}$$

Find maximum margin hyperplane

(convex optimization problem)



- Geometric margin of a hyperplane (w, b) wrt D

$$\rho(w, b) := \min_{1 \leq j \leq n} \frac{|w^T x_j + b|}{\|w\|}$$

(distance of feature vector which is closest to the hyperplane)

- Linearly Separable training set D :

$$\exists (w, b) \text{ st. } \begin{array}{l} \text{exists} \\ y_j (w^T x_j + b) > 0 \\ \forall j = 1, \dots, n \end{array}$$

$$\Rightarrow \boxed{\rho(w, b) > 0}$$

Canonical Parameterization of hyperplane (w, b) wrt lin. sep. D :

If D is lin. sep. by (w, b) and $\rho > 0$, then the canonical parameterization $(w_{\text{canon}}, b_{\text{canon}})$ is such that

$$\min_{1 \leq j \leq n} |w_{\text{canon}}^T x_j + b_{\text{canon}}| = 1$$

for canonical parameterization : $\rho(w, b) = \frac{1}{\|w\|}$

$$(w_{\text{svm}}, b_{\text{svm}}) = \underset{(w, b)}{\operatorname{argmax}} \quad \rho(w, b)$$

- 1) Perfect lin. sep. of D
- 2) Canonical rep.

$$(w_{\text{svm}}, b_{\text{svm}}) = \underset{(w, b)}{\operatorname{argmax}} \quad \frac{1}{\|w\|}$$

$$\begin{aligned} 1) \quad & y_j (w^T x_j + b) > 0 \quad \forall j \\ 2) \quad & \min_{1 \leq j \leq n} |w^T x_j + b| = 1 \end{aligned}$$

$$\rho_{\text{opt}} = \rho(w_{\text{svm}}, b_{\text{svm}})$$

$$= \max_{(w, b)} \frac{1}{\|w\|}$$

1) ...

2) ...

$$\rho_{\text{opt}} = \max_{(w, b)} \frac{1}{\|w\|} = \max_{(w, b)} \frac{1}{\|w\|}$$

$$1) \quad y_j (w^T x_j + b) > 0 \quad \forall j$$

$$2) \quad \min_j |w^T x_j + b| = 1$$

$$y_j (w^T x_j + b) \geq 1 \quad \forall j$$

$$2) \Rightarrow |w^T x_j + b| \geq 1$$

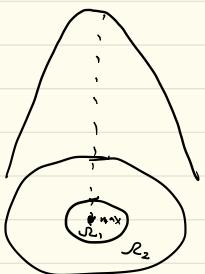
$$= y_j (w^T x_j + b) \geq 1$$

$$(\omega_*, b_*)$$

$$\frac{1}{\|\omega_*\|}$$

$$\min_{1 \leq j \leq 1} |w_*^T x_j + b_*| > 1$$

$$\text{take } (\omega_{\text{svm}}, b_{\text{svm}}) \subset \frac{(\omega_*, b_*)}{r}$$



$$\Rightarrow (\omega_{\text{svm}}, b_{\text{svm}}) = \underset{(\omega, b)}{\operatorname{argmax}} \frac{1}{\|\omega\|}$$

$\forall j, y_j (w^T x_j + b) \geq 1$

$$(\omega_{\text{svm}}, b_{\text{svm}}) = \underset{(\omega, b)}{\operatorname{argmin}} \frac{1}{2} \|\omega\|^2$$

$\forall j, y_j (w^T x_j + b) \geq 1$

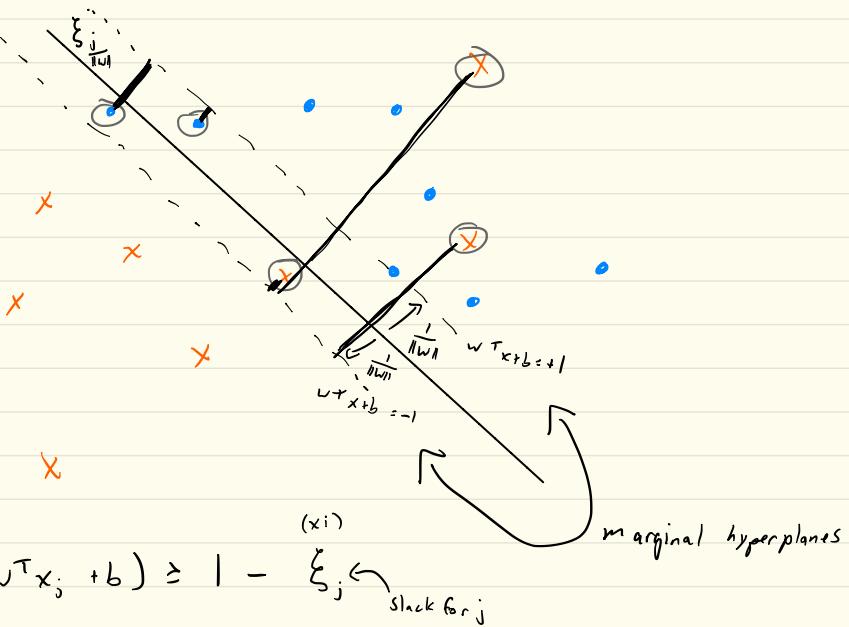
$$= \underset{(\omega, b)}{\operatorname{argmin}} f(\omega, b)$$

$\forall j, g_j(\omega, b) \leq 0$

$$f(\omega, b) = \frac{1}{2} \|\omega\|^2$$

$$g_j(\omega, b) = 1 - y_j (w^T x_j + b)$$

Non-separable



Objective function of SVM:

$$(w_{SVM}, b_{SVM}) = \underset{(w,b)}{\operatorname{arg\,min}} \left[\frac{1}{2} \|w\|^2 + \text{Penalty}(\xi_1, \dots, \xi_n) \right]$$

$$\forall j \quad y_j (w^T x_j + b) \geq 1 - \xi_j, \quad \xi_j \geq 0$$

Penalty Examples

$$1) \quad C \left(\sum_{j=1}^n \xi_j \right), \quad C > 0$$

$$2) \quad C \sum_{j=1}^n (\xi_j)^p, \quad C > 0, \quad p > 0$$

$$3) \quad C_+ \sum_{j: y_j = +1} \xi_j + C_- \sum_{j: y_j = -1} \xi_j$$

$$w_{SVM} = \sum_{j=1}^n \alpha_j^* y_j x_j, \quad 0 \leq \alpha_j^* \leq C_j \quad (\forall j)$$

$$b_{SVM} = y_j - w_{SVM}^T x_j \quad \text{for any } j: 0 < \alpha_j^* < C_j$$

$$\text{Typically } n_+ C_+ = n_- C_- = C$$

$$\Rightarrow C_+ = \frac{C}{n_+}, \quad C_- = \frac{C}{n_-}$$

n_+ = # pos samples, n_- = # neg samples

General case :

$$\text{Penalty} = \sum_{j=1}^n c_j \xi_j, \quad c_j > 0 \quad \forall j$$

Primal optimization problem

$$(w_{\text{svm}}, b_{\text{svm}}) = \underset{(w, b)}{\operatorname{argmin}} \left[\frac{1}{2} \|w\|^2 + \sum_{j=1}^n c_j \xi_j \right]$$

$$H_j \left[\begin{array}{l} y_j (w^T x_j + b) \geq 1 - \xi_j \\ \xi_j \geq 0 \end{array} \right]$$

$\xi_j \rightarrow$ "slack" variable

$c_j > 0 \rightarrow$ trade off between

maximizing soft margin $1/\|w\|$ and

minimizing degree of ideal constraint violation (slack)

Solution

convex optimization problem : quadratic program (QP)

unique global minimizer given by:

$$w_{\text{svm}} = \sum_{j=1}^n \alpha_j^* y_j x_j$$

satisfying the following complementary

slackness conditions

$$1) \quad H_j \quad \alpha_j^* \left[y_j (w_{\text{svm}}^T x_j + b_{\text{svm}}) - 1 + \xi_j^* \right] = 0$$

$$2) \quad (c_j - \alpha_j^*) \xi_j^* = 0 \quad \forall j$$

$$b_{\text{svm}} = y_j - w_{\text{svm}}^T x_j \quad \underbrace{\text{for any } j: 0 < \alpha_j^* < c_j}$$

$$\& \underline{\alpha^*} = (\alpha_1^*, \dots, \alpha_n^*)^T$$

is the solution to the following DUAL optimization problem

Solution

• convex optimization problem: quadratic program (QP)

• unique global minimizer given by:

$$w_{\text{SVM}} = \sum_{j=1}^n \alpha_j^* y_j x_j$$

Satisfying the following Complementary Slackness conditions

$$1) \quad \forall j \quad \alpha_j^* \left[y_j (w_{\text{SVM}}^T x_j + b_{\text{SVM}}) \right] = 0$$

$$2) \quad (\gamma_j - \alpha_j^*) \xi_j^* = 0 \quad \forall j$$

$$b_{\text{SVM}} = y_j - w_{\text{SVM}}^T x_j \quad \underbrace{\text{for any } j: 0 < \alpha_j^* < \gamma_j}$$

$$\& \quad \underline{\alpha}^* = (\alpha_1^*, \dots, \alpha_n^*)^T$$

is the solution to the following DUAL optimization problem

$$\underline{\alpha}^* = \arg \max_{\underline{\alpha}:} \left[\sum_{j=1}^n \alpha_j - \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \alpha_j \alpha_k y_j y_k (x_j^T x_k) \right]$$

$$\begin{array}{l} \uparrow \\ \forall j, \quad 0 \leq \alpha_j \leq \gamma_j \\ \sum_{j=1}^n \alpha_j y_j = 0 \end{array}$$

Can be solved for using sequential minimal optimization (SMO) algorithm

$$O(n^2) \quad \gamma \in [2, 3]$$

$\{x_j : \alpha_j > 0\} \rightarrow$ Support Vectors

$\{x_j : \xi_j > 0\} \rightarrow$ Outliers

All outliers are support vectors

\rightarrow from complementary slackness condition #2: $\xi_j^* > 0 \Rightarrow \alpha_j^* = \gamma_j > 0$

Support vectors $\Rightarrow \alpha_j^* \neq 0$

Then by Comp. Slackness cond #1

$$y_j (\omega_{\text{SVM}}^\top x_j + b_{\text{SVM}}) - 1 + \xi_j^* = 0$$



$$\xi_j^* > 0$$

$$\text{If } \xi_j^* = 0$$

(not outlier)

\Rightarrow outlier

$$y_j (\omega_{\text{SVM}}^\top x_j + b_{\text{SVM}}) = +1$$

Unconstrained form of the SVM Primal Optimization Problem

$$\forall j \left\{ y_j (\omega^\top x_j + b) \geq 1 - \xi_j, \quad \xi_j \geq 0 \right.$$

$$\left. \forall j \left[\xi_j \geq 1 - y_j (\omega^\top x_j + b), \quad \xi_j \geq 0 \right] \right]$$

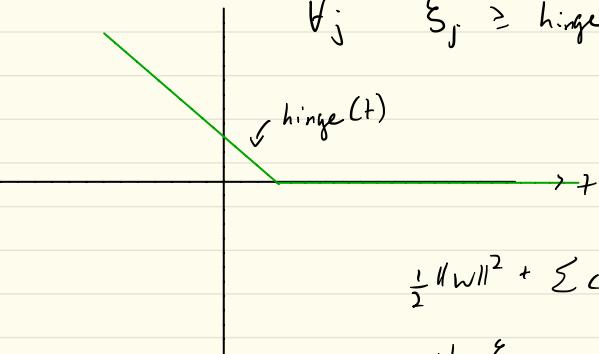
$$= \forall j \left[\xi_j \geq \max(0, 1 - y_j (\omega^\top x_j + b)) \right]$$

$$\text{hinge}(y_j(w^\top x_j + b))$$

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

$$(w_{\text{SVM}}, b_{\text{SVM}}) = \underset{(w, b)}{\operatorname{arg\,min}} \left[\frac{1}{2} \|w\|^2 + \sum_{i=1}^n c_i \xi_i \right]$$

$$y_j \xi_j \geq \text{hinge}(y_j(w^\top x_j + b))$$



$$\frac{1}{2} \|w\|^2 + \sum c_i \xi_i > \frac{1}{2} \|w\|^2 + \sum_{i=1}^n c_i \text{hinge}(y_i(w^\top x_i + b))$$

$$y_j \xi_j \geq \text{hinge}(y_j(w^\top x_j + b))$$

Ex.

$$\text{cost} = \frac{1}{2} \|w\|^2 + C_1 \text{hinge}(1) + C_2 \text{hinge}(2)$$

$$C_1 \rightarrow 2 C_1$$

ω_{svm} is unique

Proof:

$$\begin{pmatrix} \omega \\ b \end{pmatrix} \in \min_c [f(\omega) + g(b)]$$

\uparrow convex set \downarrow strictly convex \downarrow convex

b_{svm} is unique if D is linearly separable

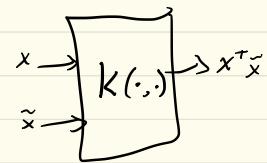
If D is not lin. sep. there may be multiple (even infinitely many) solutions for b_{svm}



$$\omega_{\text{svm}} = \begin{pmatrix} c \\ 0 \end{pmatrix}$$

$$h_{\text{SVM}}(x) = \text{Sign}(\omega_{\text{SVM}}^T x + b_{\text{SVM}})$$

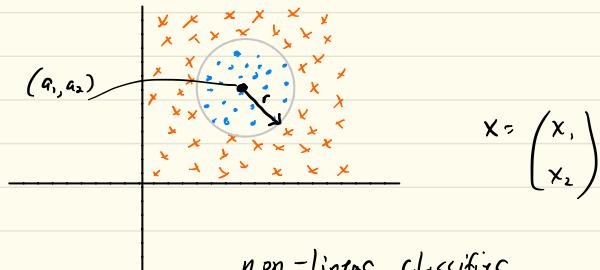
$$= \text{Sign}\left(\sum_{i=1} \alpha_i^* y_i (x_i^T x) + b_{\text{SVM}}\right)$$



Kernel Trick Motivation

- \mathcal{D} may not be linearly separable, but may be so "non-linearly"

Example



non-linear classifier

$$(x_1 - a_1)^2 + (x_2 - a_2)^2 \stackrel{+1}{\leq} r^2 \stackrel{-1}{\geq}$$

$$a_1^2 + x_2^2 - 2a_1 x_1 - 2a_2 x_2 + \underbrace{(a_1^2 + a_2^2 - r^2)}_{-1} \stackrel{+1}{\leq} 0 \stackrel{-1}{\geq}$$

$$\underbrace{\begin{pmatrix} 1 & 1 & -2a_1 & -2a_2 \end{pmatrix}}_{w^T} \begin{pmatrix} x_1^2 \\ x_2^2 \\ x_1 \\ x_2 \end{pmatrix} = b$$

$$\Phi\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right)$$

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$D_{\text{new}} = \{(\Phi(x_1), y_1), \dots, (\Phi(x_n), y_n)\}$$

Idea replace $x \in \mathbb{R}^d$ with $\phi(x) \in \mathbb{R}^{d'}$, $d' > d$, to make D more linearly separable

- But this also increases computational complexity (dot prods now in $\mathbb{R}^{d'}$)
- We also need to construct $\phi(\cdot)$

Solution : Positive Definite Symmetric (PDS) kernels

- these IMPLICITLY define an inner product in an IMPLICIT high-dim space (could even be ∞ -dimensional)

PDS kernel definition

$K(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is PDS if

1) Symmetric : $\forall x, \tilde{x} \in \mathbb{R}^d$
 $K(x, \tilde{x}) = K(\tilde{x}, x)$

2) $\forall n$ & all $x_1, \dots, x_n \in \mathbb{R}^d$
the $n \times n$ matrix

$$K_n = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix}$$

is positive semi-definite

i.e. $\underline{\zeta} \subseteq \in \mathbb{R}^n$, $\underline{\zeta}^T K_n \underline{\zeta} \geq 0$

Result (Mercer's theorem in functional analysis)

If $K(\cdot, \cdot)$ is PDS then

$$\exists \quad \phi : \mathbb{R}^d \rightarrow \mathcal{V}$$

Vector space of
suitable dimensions
(could be ∞)

Hilbert
space

such that $\forall x, \tilde{x} \in \mathbb{R}^d$

$$K(x, \tilde{x}) = \langle \phi(x), \phi(\tilde{x}) \rangle$$

(suitable inner product in \mathcal{V})

Examples of PDS kernels

1) Polynomial

$$K(x, \tilde{x}) := (x^T \tilde{x} + c)^l$$

$c > 0$ offset, $l = \text{positive integer}$

Special case $c=0, l=1$

$$K(x, \tilde{x}) = x^T \tilde{x}$$
$$\equiv \langle \phi(x), \phi(\tilde{x}) \rangle$$

$\phi(x) = x$ Identity Transformation

- Gaussian Kernel / Radial Basis Function (RBF) Kernel

$$K(x, \tilde{x}) := e^{-\frac{\|x - \tilde{x}\|^2}{2\sigma^2}}, \quad \sigma > 0$$

- Normalized RBF

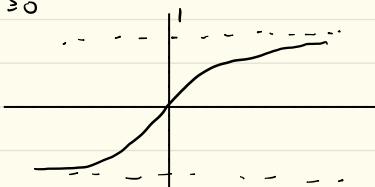
$$K(x, \tilde{x}) := e^{\frac{x^T \tilde{x}}{\sigma^2}}, \quad \sigma > 0$$

- Sigmoid kernel

$$K(x, \tilde{x}) := \tanh(a(x^T \tilde{x}) + b)$$

$a > 0, \quad b \geq 0$

$$\tanh(t) = \left(\frac{e^t - e^{-t}}{e^t + e^{-t}} \right)$$



Properties (PDS kernel algebra):

Let $K_1(\cdot, \cdot)$ & $K_2(\cdot, \cdot)$ be two PDS kernels

then the following new kernels derived from K_1 & K_2 are also PDS

1) sum : $K_{\text{sum}}(x, \tilde{x}) := K_1(x, \tilde{x}) + K_2(x, \tilde{x})$

2) positive scaling : $\alpha K_1(x, \tilde{x}), \quad \forall \alpha > 0$

3) product : $K_{\text{prod}}(x, \tilde{x}) := (K_1(x, \tilde{x}))(K_2(x, \tilde{x}))$

$$4) \underline{\text{exponentiation}} : K_{\exp}(x, \tilde{x}) := e^{k_1(x, \tilde{x})}$$

$$5) \underline{\text{Power-series}} : K_{\text{Pws}}(x, \tilde{x}) = \sum_{l=0}^{\infty} a_l (k, (x, \tilde{x}))^l$$

Multi-class SVMs

↪ Aggregate / Combined algorithms

↪ One-vs-all (OVA)

↪ One-vs-one (OVO)

↪ Error-correction-classification (ECC)

One-vs-all (OVA)

Idea: m classes : 1, 2, ..., m

Learn m Binary classifiers

For $k = 1, \dots, m$

$$h_k(x) = \begin{cases} +1 & \Rightarrow x \text{ classified as class } k \\ -1 & \Rightarrow x \text{ classified as } \underline{\text{NOT}} \ k \end{cases}$$

For k^{th} classifier: relabel training data as follows

$$y_i(k) = \begin{cases} +1 & \text{if } y_i = k \\ -1 & \text{otherwise} \end{cases}$$

Final Decision

Typically $h_k(x) = \underbrace{\text{sign}(f_k(x))}_{\text{Some scoring function}}$

$h_{\text{OVA}}(x) = \underset{1 \leq k \leq m}{\operatorname{argmax}} f_k(x)$
--

issues

1) Unbalanced data sets

→ If all classes have approx. equal examples then $\forall k$,
 the # positive samples for k^{th} classifier = $\frac{n}{m}$
 vs. $\frac{(m-1)}{m}n$ for negative

2) Uncalibrated score functions $f_1(x) \dots f_m(x)$

→ not a probability (confidence)

So range of f_1 may be different from f_2 etc.

One-vs.-one (OVO)

$$m \text{ classes} \Rightarrow \binom{m}{2} = \frac{m(m-1)}{2} \text{ pairs of classes}$$

- Learn $\binom{m}{2}$ binary classifiers.

$h_{\{k, \tilde{k}\}}(x)$ one for every unordered pair $\{k, \tilde{k}\}$, $k \neq \tilde{k}$

$$h_{\{k, \tilde{k}\}}(x) = \begin{cases} k & \text{if classified as } k \\ \tilde{k} & \text{otherwise} \end{cases}$$

$$h_{\text{ovo}}(x) = \underset{1 \leq k \leq m}{\arg \max} \underbrace{\left| \left\{ k \neq \tilde{k} : h_{\{k, \tilde{k}\}}(x) = k \right\} \right|}_{\text{\# of classes against which class } k \text{ 'wins'}}$$

↑
class with the largest
number of wins

(+) → No calibration problem
Combining decisions not uncalib. scores

(+/-?) → Computational Complexity
- training
- testing

(-) → Fewer total # points for training any individual class
 $\left(\frac{2n}{m}\right)$

(+) → Datasets are balanced for each classifier

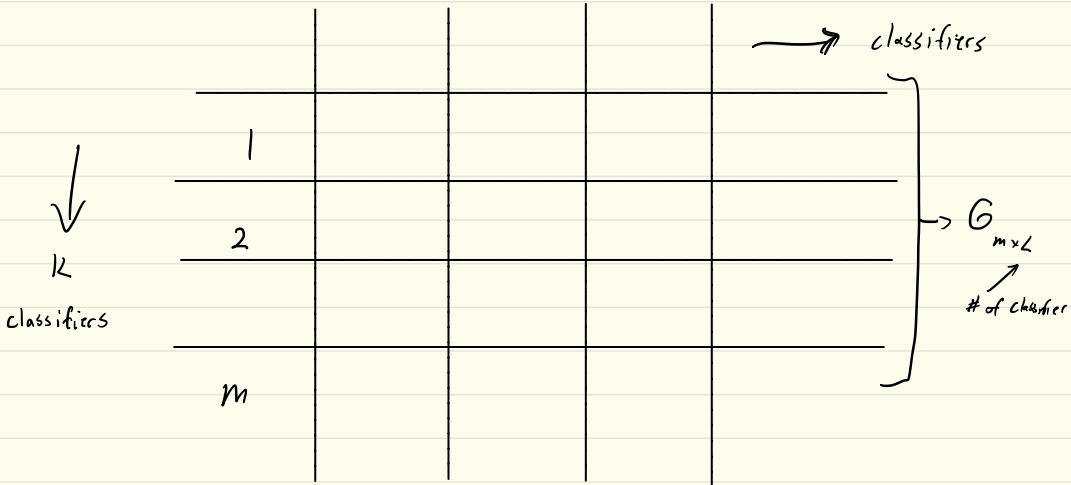
Computational Cost (rough) analysis

Assume $\approx \frac{n}{m}$ pts. in each class

$O(n^8)$ complexity for training a single binary classifier in the family on a dataset of size n

$C_{\text{test}} = \text{Cost of testing with 1 binary}$

	training	testing
OVA	$O(m n^8)$	$O(m C_{\text{test}})$
OVO	$O(m^2 \left(\frac{n}{m}\right)^8)$	$O(m^2 C_{\text{test}})$



Multi-class classification (not specific to SVM)

Combined / Aggregated Methods

- Family of binary classifiers
- trained on different labelings of training set
- scores/ decisions are combined/ aggregated to arrive at final multiclass label

- | |
|---|
| <ul style="list-style-type: none"> - OVA - OVO - ECC |
|---|

Uncombined multi-class SVM

2 class

$$(\underline{w}_{\text{svm}}, b_{\text{svm}}) = \underset{(\underline{w}, b, \xi)}{\operatorname{argmin}} \left[\frac{1}{2} \|\underline{w}\|^2 + \sum_{j=1}^n c_j \xi_j \right]$$

$$\forall j : y_j (\underline{w}^\top x_j + b) \geq 1 - \xi_j$$

$$\xi \geq 0$$

$$h_{\text{svm}}(x) = \operatorname{sign} (\underline{w}_{\text{svm}}^\top x + b_{\text{svm}})$$

$$h_{\text{multiclass SVM}}(x) = \underset{1 \leq k \leq m}{\operatorname{argmax}} \left(\underline{w}_{k,\text{svm}}^\top x + b_{k,\text{svm}} \right)$$

$$(\underline{w}_1, \text{svm}, \underline{w}_2, \text{svm}, \dots, \underline{w}_m, \text{svm})$$

$$(b_1, \text{svm}, b_2, \text{svm}, \dots, b_m, \text{svm})$$

$$= \operatorname{argmin} \left[\frac{1}{2} \sum_{k=1}^m \|\underline{w}_k\|^2 + \sum_{j=1}^n c_j \xi_j \right]$$

$$\forall j, \forall l \neq y_j, \xi_j \geq 0,$$

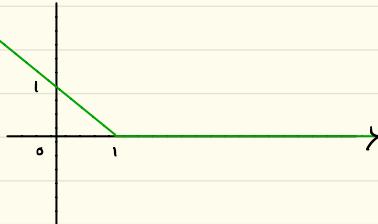
$$(\underline{w}_{y_j}^\top x + b_{y_j}) \geq (\underline{w}_l^\top x + b_l) + 1 - \xi_j$$

SVMs do not arise from either Generative / Discriminative probabilistic models

Unconstrained Form of the SVM optimization problem

$$(\omega_{\text{svm}}, b_{\text{svm}}) = \underset{(\omega, b)}{\operatorname{argmin}} \left[\frac{1}{2} \|\omega\|^2 + \sum_{j=1}^n c_j \operatorname{hinge}(y_j (\omega^\top x_j + b)) \right]$$

$$\operatorname{hinge}(t) = \max(0, 1-t)$$



Assume (only for simplicity) that $c_j = c > 0 \ \forall j$

then minimizing becomes

$\sum_{j=1}^n c \cdot \operatorname{hinge}(y_j (\omega^\top x_j + b))$ }	$+ \frac{1}{2} \ \omega\ ^2$ } $- \ln \pi(\omega)$
Gen. : $-\ln p(x_i, y_i \overset{\theta}{\omega}, b)$	$\pi(\omega) = \frac{e^{-\ \omega\ ^2}}{\sqrt{?}}$
Disc. : $-\ln p(y_i x_i, \omega, b)$	

Gen?

$$p(x, y | w, b) = \frac{1}{2} e^{-\text{chinge}(y(w^T x + b))}$$

Require

$$\sum_{y=\pm 1} \int_{x \in \mathbb{R}^d} p(x, y | w, b) dx = 1$$

$$\Rightarrow \frac{1}{2} \int_x \left[e^{-\text{chinge}(+(w^T x + b))} + e^{-\text{chinge}(-(w^T x + b))} \right] dx = 1$$

$Z =$