# Project Presentations

**Tuesday, 3/1**
1. Group 1 (**Hussain**, Tanvir; **Lewis**, Cameron; **Villamar**, Sandra)
2. Group 2 (**Dong**, Meng; **Long**, Jianzhi; **Wen**, Bo; **Zhang**, Haochen)
3. Group 3 (**Chen**, Yuzhao; **Li**, Zonghuan; **Song**, Yuze; **Yan**, Ge)
4. Group 4 (**Li**, Jiayuan; **Xiao**, Nan; **Yu**, Nancy; **Zhou**, Pei)
5. Group 5 (**Li**, Zheng; **Tao**, Jianyu; **Yang**, Fengqi)
6. Group 6 (**Bian**, Xintong; **Jiang**, Yufan; **Wu**, Qiyao)
7. Group 7 (**Chen**, Yongxing; **Yao**, Yanzhi; **Zhang**, Canwei)
8. Group 8 (**Nukala**, Kishore; **Pulleti**, Sai; **Vaidyula**, Srikar)

**Thursday, 3/3**
1. Group 9 (**Baluja**, Michael; **Cao**, Fangning; **Huff**, Mikael; **Shen**, Xuyang)
2. Group 10 (**Arun**, Aditya; **Long**, Heyang; **Peng**, Haonan)
3. Group 11 (**Cowin**, Samuel; **Liao**, Albert; **Mandadi**, Sumega)
4. Group 12 (**Jia**, Yichen; **Jiang**, Zhiyun; **Li**, Zhuofan)
5. Group 13 (**Dandu**, Murali; **Daru**, Srinivas; **Pamidi**, Sri)
6. Group 14 (**He**, Bolin; **Huang**, Yen-Ting; **Wang**, Shi; **Wang**, Tzu-Kao)
7. Group 15 (**Chen**, Luobin; **Feng**, Ruining; **Wu**, Ximei; **Xu**, Haoran)

**Tuesday, 3/8**
1. Group 16 (**Chen**, Rex; **Liang**, Youwei; **Zheng**, Xinran)
2. Group 17 (**Aguilar**, Matthew; **Millhiser**, Jacob; **O'Boyle**, John; **Sharpless**, Will)
3. Group 18 (**Wang**, Haoyu; **Wang**, Jiawei; **Zhang**, Yuwei)
4. Group 19 (**Chen**, Yinbo; **Di**, Zonglin; **Mu**, Jiteng)
5. Group 20 (**Chowdhury**, Debalina; **He**, Scott; **Ye**, Yiheng)
6. Group 21 (**Lin**, Wei-Ru; **Ru**, Liyang; **Zhang**, Shaohua)
7. Group 22 (**Bhavsar**, Shivad; **Blazej**, Christopher; **Bu**, Yinyan; **Liu**, Haozhe)

**Thursday, 3/10**
1. Group 23 (**Chen**, Claire; **Hsieh**, Chia-Wei; **Lin**, Jui-Yu; **Tsai**, Ya-Chen)
2. Group 24 (**Cheng**, Yu; **Yu**, Zhaowei; **Zaidi**, Ali)
3. Group 25 (**Assadi**, Parsa; **Brugere**, Tristan; **Pathak**, Nikhil; **Zou**, Yuxin)
4. Group 28 (**Candassamy**, Gokulakrishnan; **Dixit**, Rajeev; **Huang**, Joyce)
5. Group 27 (**Kok**, Hong; **Wang**, Jacky; **Yan**, Yijia; **Yuan**, Zhouyuan )
6. Group 28 (**Luan**, Zeting; **Yang**, Zheng)
7. Group 29 (**Cuawenberghs**, Kalyani; **Mojtahed**, Hamed)

Each presentation will be allocated 9 minutes (pts will be deducted if you go over 9 minutes)

The presentation slides of ALL GROUPS are due by Monday, 2/28 @ 11:59 pm

Email me the file (mvasconcelos@eng.ucsd.edu) and name the file GroupX.pdf, where X is your group number (see previous slide). Use Group X Presentation as the subject of your email and cc to all members.

The presentation should discuss the **problem** that you are trying to solve, the **data that you are using**, the **proposed solution**(s), and the **results** that you have so far (they can later be UPDATED IN THE PROJECT PAPER).

# ECE 271B – Winter 2022

# The Soft−Margin Support Vector Machine

Disclaimer:
This class will be recorded
and made available to students asynchronously.

Manuela Vasconcelos

ECE Department, UCSD

# The Support Vector Machine

▶ the SVM is the **classifier** that <u>maximizes the margin</u> under the constraints

$$\min_{\mathbf{w}, b} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1, \forall i$$

▶ <u>no dual gap</u>, and the **dual problem** is

$$\max_{\boldsymbol{\alpha} \geq 0}\left\{-\frac{1}{2}\sum_{ij} \alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j + \sum_i \alpha_i\right\} \quad \text{subject to} \quad \sum_i \alpha_i y_i = 0$$

once this is solved, the vector

$$\mathbf{w}^* = \sum_i \alpha_i^* y_i \mathbf{x}_i$$

is the <u>normal</u> to the **maximum margin plane**

▶ **note:** the dual solution does <u>not</u> determine the optimal $b^*$

2

# Support Vectors

▶ from the **KKT conditions**, a <span style="color:magenta">active</span> (<span style="color:green">inactive</span>) constraint has <span style="color:magenta">non−zero</span> (<span style="color:green">zero</span>) Lagrange multiplier $\alpha_i$

▶ that is

$$\alpha_i > 0 \ \ \text{iff} \ \ y_i\big(\mathbf{w}^{*T}\mathbf{x}_i + b^*\big) = 1$$

▶ hence $\boxed{\alpha_i > 0}$ only for points

$$\boxed{\big|\mathbf{w}^{*T}\mathbf{x}_i + b^*\big| = 1}$$

which are those **that lie at a distance** <u>equal</u> to the **margin**

▶ these "support" the **hyperplane** and are called <u>**support vectors**</u>

▶ the **decision rule** is

$$f(\mathbf{x}) = \text{sgn}\left[\sum_{i \in SV} y_i \alpha_i^* \mathbf{x}_i^T \mathbf{x} + b^*\right]$$

$$SV = \{i \mid \alpha_i^* > 0\}$$

and the **remaining points are** <u>**irrelevant**</u>!

$\alpha_i = 0$

$\gamma$

$\alpha_i > 0$

$\alpha_i = 0$

# Hard−Margin SVM

▶ **SVM training**

1) solve the **optimization problem**

$$\max_{\boldsymbol{\alpha} \geq 0}\left\{-\frac{1}{2}\sum_{ij}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \alpha_i\right\} \quad \text{subject to} \quad \sum_i \alpha_i y_i = 0$$

2) then **compute**

$$\mathbf{w}^* = \sum_{i \in SV} \alpha_i^* y_i \mathbf{x}_i$$

$$b^* = -\frac{1}{2}\sum_{i \in SV} y_i \alpha_i^* \left(\mathbf{x}_i^T \mathbf{x}^+ + \mathbf{x}_i^T \mathbf{x}^-\right)$$

▶ **decision function**

$$f(\mathbf{x}) = \text{sgn}\left[\sum_{i \in SV} y_i \alpha_i^* \mathbf{x}_i^T \mathbf{x} + b^*\right]$$
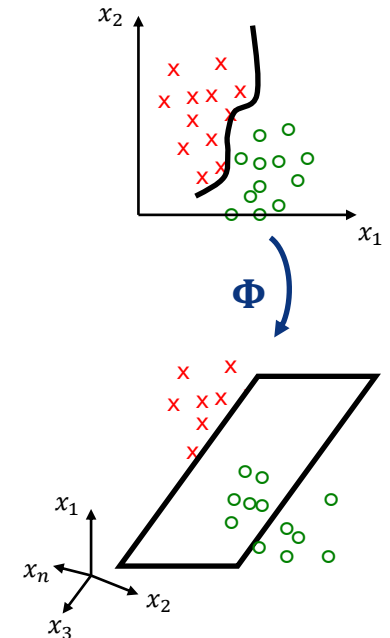
4

# SVM: Kernelization

$$\max_{\boldsymbol{\alpha} \geq 0}\left\{-\frac{1}{2}\sum_{ij}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j + \sum_i \alpha_i\right\} \quad \text{subject to} \quad \sum_i \alpha_i y_i = 0$$

$$b^* = -\frac{1}{2}\sum_{i \in SV} y_i \alpha_i^*\left(\mathbf{x}_i^T\mathbf{x}^+ + \mathbf{x}_i^T\mathbf{x}^-\right)$$

- note that **all** equations <u>depend</u> **only** on $\boxed{\mathbf{x}_i^T \mathbf{x}_j}$

$$f(\mathbf{x}) = \text{sgn}\left[\sum_{i \in SV} y_i \alpha_i^* \mathbf{x}_i^T\mathbf{x} + b^*\right]$$

- the "kernel trick" is **trivial**: replace by $K(\mathbf{x}_i, \mathbf{x}_j)$

1) training

$$\max_{\boldsymbol{\alpha} \geq 0}\left\{-\frac{1}{2}\sum_{ij}\alpha_i\alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum_i \alpha_i\right\} \quad \text{subject to} \quad \sum_i \alpha_i y_i = 0$$

$$b^* = -\frac{1}{2}\sum_{i \in SV} y_i \alpha_i^*\left(K(\mathbf{x}_i, \mathbf{x}^+) + K(\mathbf{x}_i, \mathbf{x}^-)\right)$$

2) decision function

$$f(\mathbf{x}) = \text{sgn}\left[\sum_{i \in SV} y_i \alpha_i^* K(\mathbf{x}_i, \mathbf{x}) + b^*\right]$$

- note that we can <u>no longer</u> recover $\mathbf{w}^*$ <u>explicitly</u> without determining the feature transformation $\boldsymbol{\Phi}$, but, **luckily**, we do <u>not</u> really need $\mathbf{w}^*$, <u>only</u> the decision function

$$\mathbf{w}^* = \sum_{i \in SV} \alpha_i^* y_i \, \boldsymbol{\Phi}(\mathbf{x}_i)$$

# Input−Space Interpretation

$$f(\mathbf{x}) = \text{sgn}\left[\sum_{i \in SV} y_i \alpha_i^* K(\mathbf{x}_i, \mathbf{x}) + b^*\right]$$

▶ last class, we saw that the decision function identical to the **BDR** for

1) **class 1** with **likelihood**                    and **prior**

$$\sum_{i \in SV | y_i \geq 0} \pi_i^* K(\mathbf{x}_i, \mathbf{x})$$

$$\pi_i^* = \frac{\alpha_i^*}{\sum_{i \in SV | y_i \geq 0} \alpha_i^*}, i | y_i \geq 0$$

$$\sum_{i \in SV | y_i \geq 0} \alpha_i^* \Big/ \sum_i \alpha_i^*$$

2) **class 2** with **likelihood**                    and **prior**

$$\sum_{i \in SV | y_i < 0} \beta_i^* K(\mathbf{x}_i, \mathbf{x})$$

$$\beta_i^* = \frac{\alpha_i^*}{\sum_{i \in SV | y_i < 0} \alpha_i^*}, i | y_i < 0$$

$$\sum_{i \in SV | y_i < 0} \alpha_i^* \Big/ \sum_i \alpha_i^*$$

▶ i.e.

$$f(\mathbf{x}) = \begin{cases} 1, & \dfrac{\sum_{i \in SV | y_i \geq 0} \pi_i^* K(\mathbf{x}_i, \mathbf{x})}{\sum_{i \in SV | y_i < 0} \beta_i^* K(\mathbf{x}_i, \mathbf{x})} \geq T \\ \\ -1, & \text{otherwise} \end{cases}$$

BDR threshold

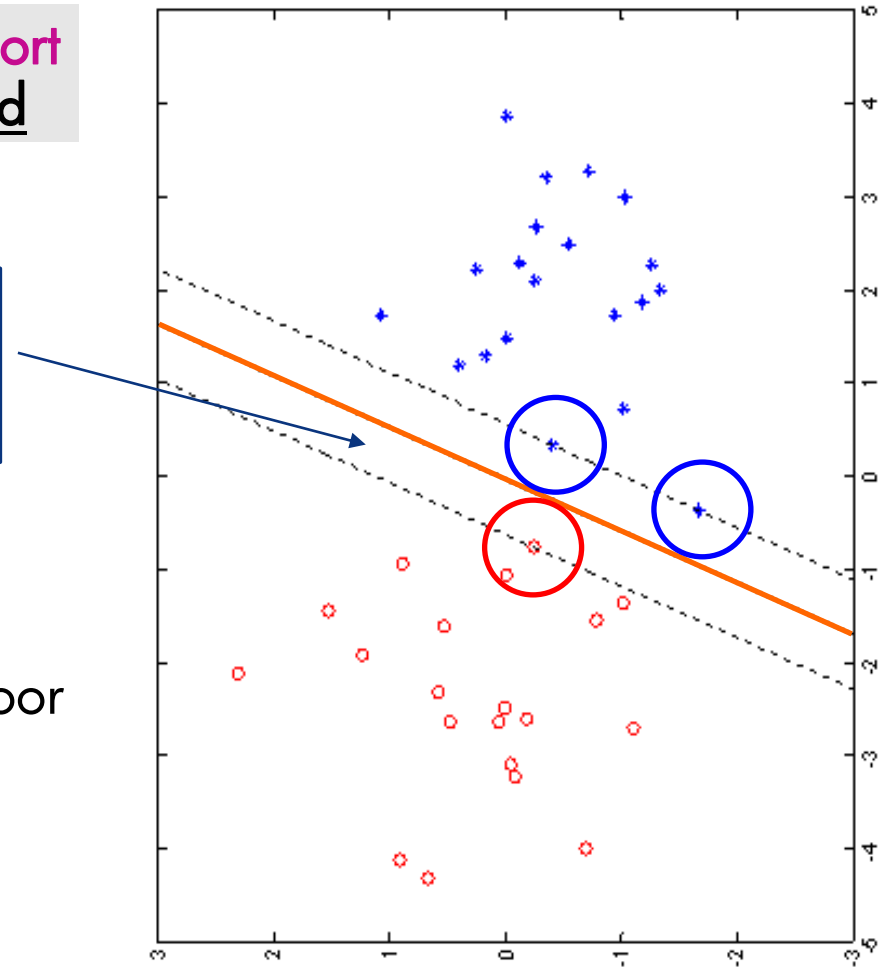▶ these **likelihood** functions are a **kernel** density estimate if $K(\cdot, \mathbf{x}_i)$ is a valid pdf

6

# Input–Space Interpretation

▶ <u>peculiar</u> kernel estimates

- <u>only</u> place kernels on the support vectors, all other points <u>ignored</u>

▶ <u>discriminant</u> density estimation
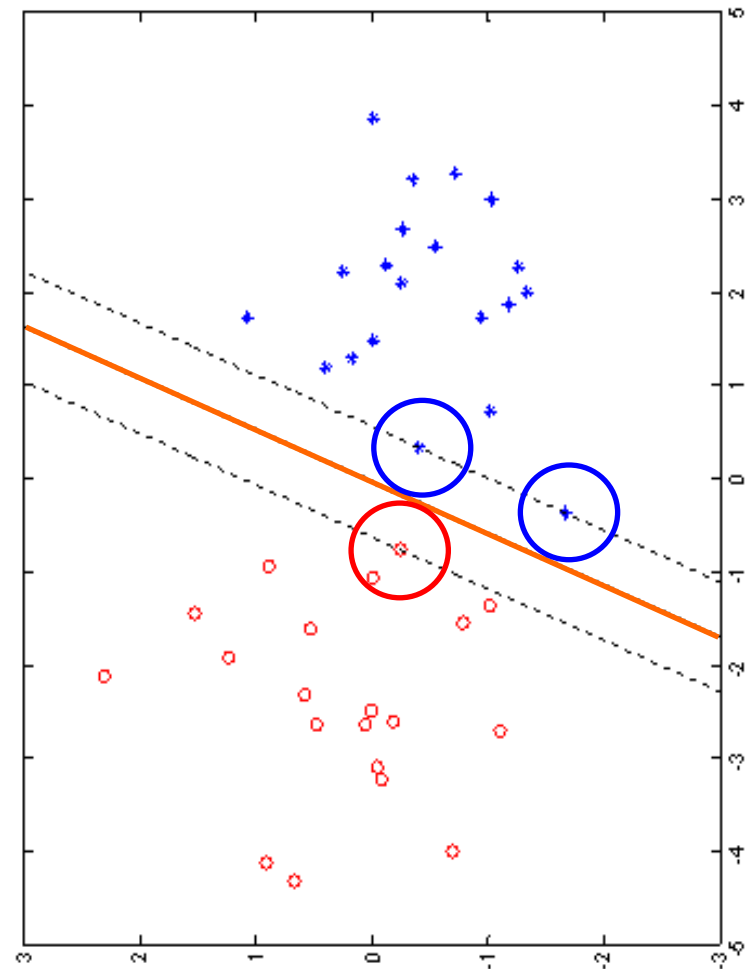
- concentrate <u>modeling power</u> where it <u>matters the most</u>, i.e. near **classification boundary**

- **smart**, since points away from the boundary are always well classified even if the density estimates in their region are poor

- the **SVM** is a highly efficient <u>combination</u> of the **BDR** with **kernel estimates**, complexity $O(|SV|)$ instead of $O(n)$
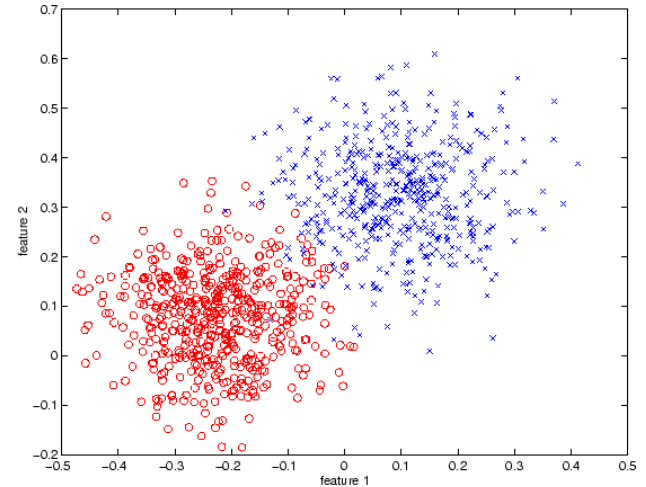
# Limitations of the SVM

▶ appealing, but also points out the <u>limitations</u> of the SVM:

- **major problem** of kernel density estimation is the choice of bandwidth

  - if too small, the estimates have too much variance

  - if too large, the estimates have too much bias

- this problem **appears** <u>again</u> in the SVM

  - <u>no</u> generic "optimal" procedure to find the kernel or its parameters

  - requires <u>trial and error</u>

  - note, however, that this is **less of a headache** since only a <u>few</u> kernels have to be evaluated



- usually, we pick an **arbitrary kernel**, e.g. Gaussian

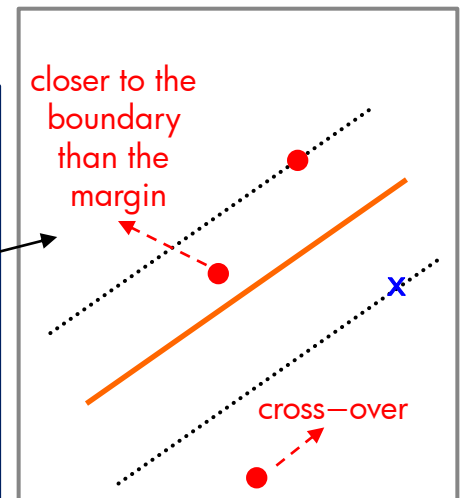- then, determine **kernel parameters**, e.g. variance, by <u>trial and error</u>

# SVM: Non−Separable Problems

▶ so far, we have **assumed** <u>linearly separable</u> classes

▶ this is <u>**rarely**</u> **the case in practice**

▶ a separable problem is "**easy**": most classifiers will do well

▶ we need to be able to **extend** the SVM to the <u>non−separable</u> case



▶ <u>basic idea</u>:

- with **class overlap**, we **cannot enforce** a margin
- but we **can** enforce a <u>soft−margin</u>
  - for <u>**most**</u> points, there is a margin
  - but then there are a <u>few outliers</u> that <u>cross−over</u> or are <u>closer</u> to the boundary than the margin



closer to the boundary than the margin

cross−over

# SVM: Soft−Margin Optimization
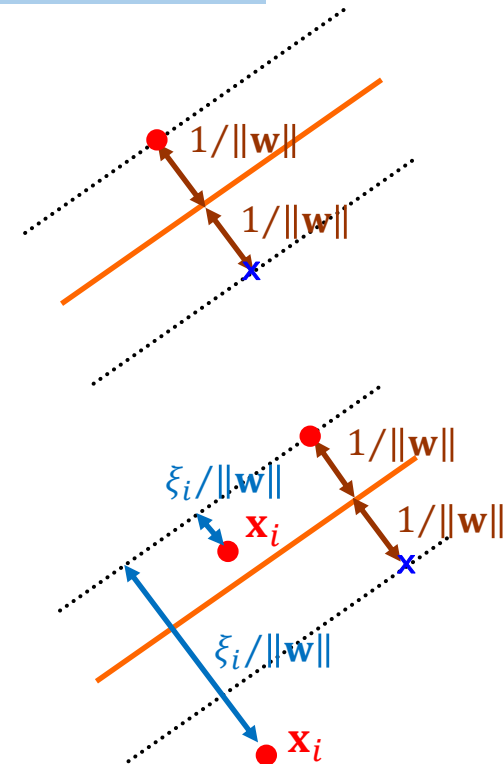
► mathematically, this can be done by introducing slack variables

► instead of solving the hard−margin problem

$$\min_{\mathbf{w},b} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1, \forall i$$
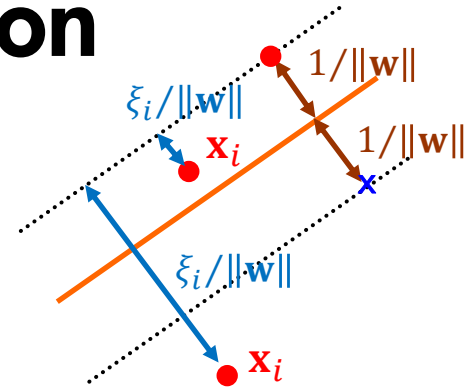
► we solve the soft−margin problem

$$\min_{\mathbf{w},\xi,b} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$$
$$\xi_i \geq 0, \forall i$$

► the $\xi_i$ are called slacks

► basically, the **same as before**, but points with $\xi_i > 0$ are allowed to violate the margin

# SVM: Soft−Margin Optimization



- note that the problem is **not** really **well defined**

  - by making $\xi_i$ arbitrarily large, any **w** will do
  - we need to penalize large $\xi_i$

- this is done by solving instead the <u>regularized optimization</u> problem

$$\min_{\mathbf{w},\boldsymbol{\xi},b} \ \|\mathbf{w}\|^2 + Cf(\boldsymbol{\xi})$$

$$\text{subject to} \quad y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$$

$$\xi_i \geq 0, \forall i$$

$Cf(\boldsymbol{\xi})$ – <u>penalty</u> or <u>regularization</u> term

$C > 0$ controls how <u>harsh</u> the penalty is

- $f(\boldsymbol{\xi})$ is **usually** a **norm**: we consider

the **1−norm**

$$f(\boldsymbol{\xi}) = \sum_i \xi_i$$

the **2−norm**

$$f(\boldsymbol{\xi}) = \sum_i \xi_i^2$$

# 2−Norm SVM

$$\min_{\mathbf{w}, \boldsymbol{\xi}, b} \ \|\mathbf{w}\|^2 + C \sum_i \xi_i^2$$

$$\text{subject to} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \forall i \qquad (**)$$

$$\xi_i \geq 0, \forall i$$

▶ **note** that

- if $\xi_i < 0$ <u>and</u> the constraint $(**)$ is satisfied,

  then $(**)$ is satisfied by $\xi_i = 0$ and the cost will be <u>smaller</u>

- hence $\xi_i < 0$ is <u>**never**</u> a solution and the positivity constraints on the $\xi_i$ are redundant

- they can therefore be dropped

12

# 2−Norm SVM

▶ this leads to

$$\min_{\mathbf{w},\boldsymbol{\xi},b} \frac{1}{2}\|\mathbf{w}\|^2 + \frac{C}{2}\sum_i \xi_i^2$$

$$\text{subject to} \quad y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$$

▶ and

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \frac{1}{2}\|\mathbf{w}\|^2 + \frac{1}{2}C\sum_i \xi_i^2 + \sum_i \alpha_i[1 - \xi_i - y_i(\mathbf{w}^T\mathbf{x}_i + b)]$$

▶ from which

$$\nabla_{\mathbf{w}}L = 0 \iff \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0 \iff \mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$\nabla_b L = 0 \iff \sum_i \alpha_i y_i = 0 \qquad \nabla_{\xi_i} L = 0 \iff C\xi_i - \alpha_i = 0$$

# 2−Norm SVM

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \tfrac{1}{2}\|\mathbf{w}\|^2 + \tfrac{1}{2} C \sum_i \xi_i^2 + \sum_i \alpha_i [1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b)]$$

▶ plugging back

$$\mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i \qquad \sum_i \alpha_i y_i = 0 \qquad \xi_i = \frac{\alpha_i}{C}$$

▶ we get the Lagrangian

$$L(\mathbf{w}^*, b, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \tfrac{1}{2}\|\mathbf{w}^*\|^2 + \tfrac{C}{2} \sum_i \left(\frac{\alpha_i}{C}\right)^2 + \sum_i \alpha_i \left[1 - \frac{\alpha_i}{C} - y_i(\mathbf{w}^{*T} \mathbf{x}_i + b)\right]$$

$$= \tfrac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \tfrac{1}{2} \sum_i \frac{\alpha_i^2}{C} + \sum_i \alpha_i - \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \underbrace{\sum_i \alpha_i y_i b}_{0}$$

$$= -\tfrac{1}{2} \left( \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \frac{1}{C} \sum_i \alpha_i^2 \right) + \sum_i \alpha_i$$

$$= -\tfrac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \left( \mathbf{x}_i^T \mathbf{x}_j + \frac{\delta_{ij}}{C} \right) + \sum_i \alpha_i \qquad \delta_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$

# Soft Dual for 2−Norm

hard−margin

$$\max_{\alpha \geq 0} \left\{ -\tfrac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \alpha_i \right\}$$

$$\text{subject to } \sum_i \alpha_i y_i = 0$$

▶ the dual problem is

$$\max_{\boldsymbol{\alpha} \geq 0} \left\{ -\tfrac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \left( \mathbf{x}_i^T \mathbf{x}_j + \frac{\delta_{ij}}{C} \right) + \sum_i \alpha_i \right\}$$

$$\text{subject to } \sum_i \alpha_i y_i = 0$$

▶ **same** as **hard−margin**, with $\frac{1}{C} \mathbf{I}$ added to kernel matrix

$$\sum_{ij} \alpha_i \alpha_j y_i y_j \left( \mathbf{x}_i^T \mathbf{x}_j + \frac{\delta_{ij}}{C} \right)$$

$$= \sum_{ij} b_i b_j K_{ij} = \mathbf{b}^T \mathbf{K} \, \mathbf{b}$$

with

$$b_i = \alpha_i y_i; \quad K_{ij} = \mathbf{x}_i^T \mathbf{x}_j + \frac{\delta_{ij}}{C}$$

▶ this:

- **increments** the eigenvalues by $1/C$, making the problem better conditioned

- for larger $C$, the extra term is **smaller** and <u>outliers</u> have a **larger** influence (**less penalty** on them, **more reliance** on data term)

# 1−Norm SVM

$$\min_{\mathbf{w}, \boldsymbol{\xi}, b} \|\mathbf{w}\|^2 + C \sum_i \xi_i \quad \text{subject to} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$$
$$\xi_i \geq 0, \forall i$$

▶ and

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{r}) =$$
$$\tfrac{1}{2}\|\mathbf{w}\|^2 + C \sum_i \xi_i + \sum_i \alpha_i [1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b)] - \sum_i r_i \xi_i$$

▶ from which

$$\nabla_{\mathbf{w}} L = 0 \iff \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0 \iff \mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$\nabla_b L = 0 \iff \sum_i \alpha_i y_i = 0 \qquad \nabla_{\xi_i} L = 0 \iff C - \alpha_i - r_i = 0$$

# 1−Norm SVM

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{r}) = \tfrac{1}{2}\|\mathbf{w}\|^2 + C\sum_i \xi_i + \sum_i \alpha_i[1 - \xi_i - y_i(\mathbf{w}^T\mathbf{x}_i + b)] - \sum_i r_i\xi_i$$

▶ plugging back

$$\mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i \qquad \sum_i \alpha_i y_i = 0 \qquad r_i = C - \alpha_i$$

▶ we get the Lagrangian

$$L(\mathbf{w}^*, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{r}) = \tfrac{1}{2}\sum_{ij} \alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j + C\sum_i \xi_i + \sum_i \alpha_i(1 - \xi_i) - \sum_{ij} \alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j$$

$$- \underbrace{\sum_i \alpha_i y_i}_{0} - \sum_i (C - \alpha_i)\xi_i$$

$$= -\tfrac{1}{2}\sum_{ij} \alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j + \sum_i \alpha_i$$

▶ this is exactly **like** the <u>hard−margin</u> case with the <u>extra</u> constraint
$\alpha_i = C - r_i, \forall i$

# Soft Dual for 1−Norm

► in **summary:**

$$\mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$\sum_i \alpha_i y_i = 0$$

$$r_i = C - \alpha_i$$

$$L(\mathbf{w}^*, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{r}) = -\frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \alpha_i$$

► <u>Recall</u>

$$\min_{\mathbf{w}, \boldsymbol{\xi}, b} \|\mathbf{w}\|^2 + C \sum_i \xi_i \quad \text{subject to} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$$
$$\xi_i \geq 0, \forall i$$

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{r}) = \frac{1}{2}\|\mathbf{w}\|^2 + C \sum_i \xi_i + \sum_i \alpha_i[1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b)] - \sum_i r_i \xi_i$$

► from the constraints

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \Leftrightarrow 1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0$$

$$\xi_i \geq 0 \Leftrightarrow -\xi_i \leq 0$$

the KKT conditions are

$$\alpha_i > 0 \Leftrightarrow 1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) = 0$$
$$\alpha_i = 0 \Leftrightarrow 1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$$

$$r_i > 0 \Leftrightarrow \xi_i = 0$$
$$r_i = 0 \Leftrightarrow \xi_i > 0$$

# Soft Dual for 1−Norm

$$L(\mathbf{w}^*, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{r}) = -\frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \alpha_i \qquad \mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i \qquad \sum_i \alpha_i y_i = 0 \qquad r_i = C - \alpha_i$$

$$(*)$$

a) $\alpha_i > 0 \Leftrightarrow 1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) = 0$      c) $r_i > 0 \Leftrightarrow \xi_i = 0$

b) $\alpha_i = 0 \Leftrightarrow 1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$      d) $r_i = 0 \Leftrightarrow \xi_i > 0$

▶ if $\boxed{\alpha_i = 0}$

- from (*), $r_i = C$ and, from c), $\xi_i = 0$
- from b), $1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$ and, since $\xi_i = 0$, we have

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) > 1, \text{ i.e. } \boxed{\mathbf{x}_i \text{ is correctly classified}}$$

▶ if $\boxed{\alpha_i > 0}$

- since $r_i$ are Lagrange multipliers, $r_i \geq 0$, (*) means that $\alpha_i \leq C$
- if $r_i > 0 \Rightarrow \boxed{\alpha_i < C}$, from c) $\xi_i = 0$

$$\text{and from a) } y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1, \text{ i.e. } \boxed{\mathbf{x}_i \text{ is on the margin}}$$

- if $r_i = 0 \Rightarrow \boxed{\alpha_i = C}$, from d) $\xi_i > 0$

$$\text{and from a) } y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 - \xi_i < 1, \text{ i.e. } \boxed{\mathbf{x}_i \text{ is an outlier}}$$

# Soft Dual for 1−Norm

$\alpha_i = 0$, $\mathbf{x}_i$ is **correctly classified**

$0 < \alpha_i < C$, $\mathbf{x}_i$ is **on the margin**

$\alpha_i = C$, $\mathbf{x}_i$ is an **outlier**

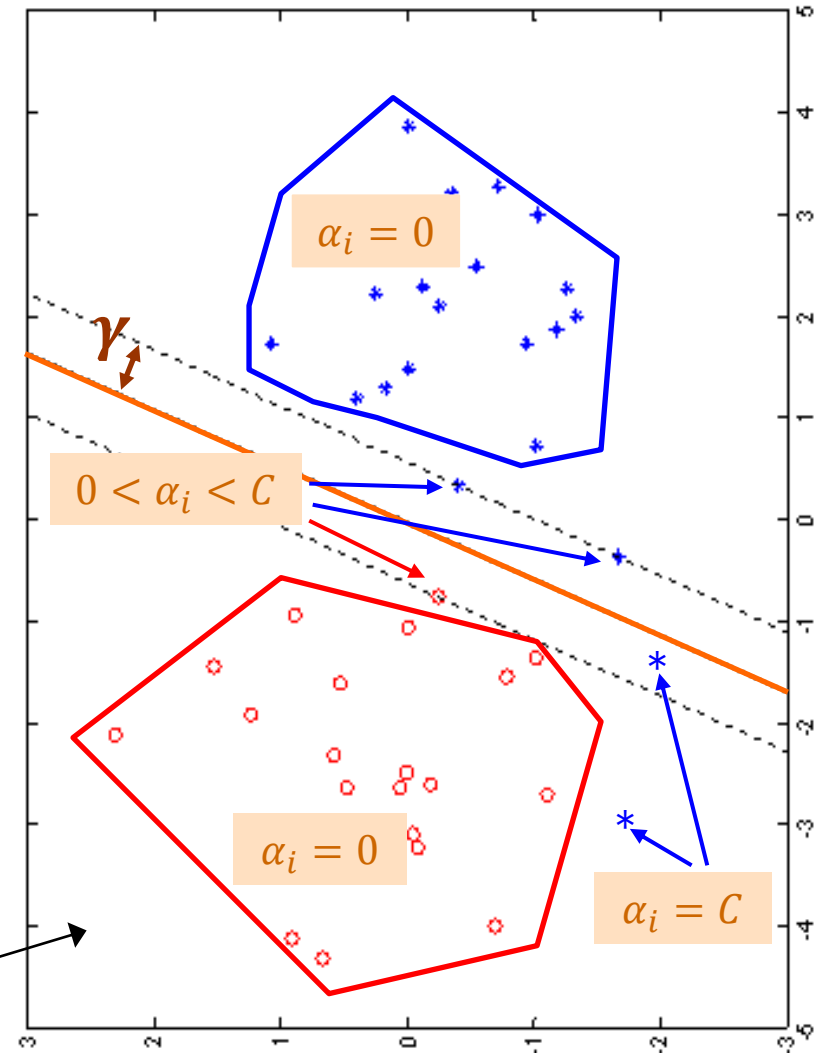▶ overall, **dual problem** is

$$\max_{\boldsymbol{\alpha} \geq 0} \left\{ -\frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \alpha_i \right\}$$

$$\text{subject to } \sum_i \alpha_i y_i = 0,$$

$$0 \leq \alpha_i \leq C$$

▶ the **only difference** with respect to the **hard−margin** case is the "**box constraint**" on the $\alpha_i$

▶ **geometrically**, we have this



$\alpha_i = 0$

$0 < \alpha_i < C$

$\gamma$

$\alpha_i = 0$

$\alpha_i = C$

# Soft Dual for 1−Norm: Support Vectors

▶ **support vectors** are the points with $\alpha_i > 0$
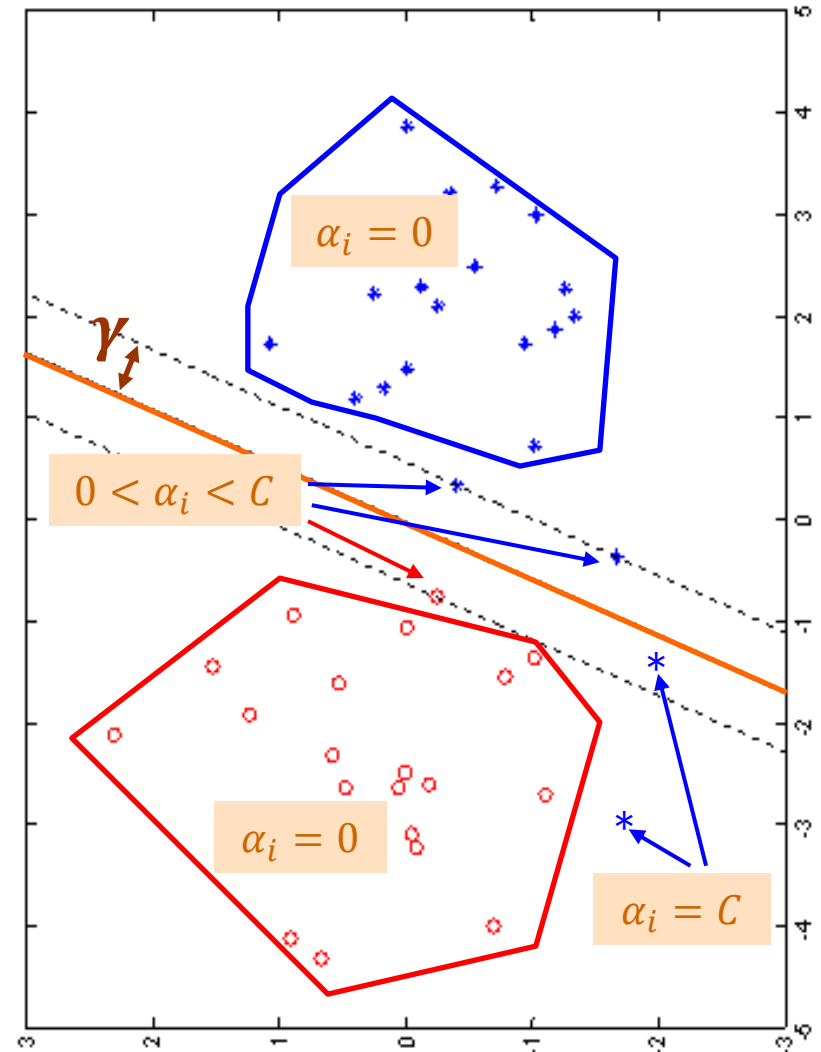
▶ as before, the **decision rule** is

$$f(\mathbf{x}) = \text{sgn}\left[\sum_{i \in SV} y_i \alpha_i^* \mathbf{x}_i^T \mathbf{x} + b^*\right]$$

where $SV = \{i \mid \alpha_i^* > 0\}$ and $b^*$ chosen such that

$$y_i g(\mathbf{x}_i) = 1, \forall \mathbf{x}_i \text{ s.t. } 0 < \alpha_i < C$$

▶ the **box constraint** on Lagrange multipliers makes **intuitive** sense:

it prevents a **single** SV outlier from having **large** impact in the decision rule



$\alpha_i = 0$

$\gamma_t$

$0 < \alpha_i < C$

$\alpha_i = 0$

$\alpha_i = C$

# Soft−Margin SVM

$$\min_{\mathbf{w},\xi,b} \quad \|\mathbf{w}\|^2 + Cf(\xi)$$
$$\text{subject to} \quad y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$$
$$\xi_i \geq 0, \forall i$$

► note that $C$ **controls** the <u>importance</u> of outliers

- larger $C$ implies that <u>**more**</u> emphasis is given to **minimizing the number of outliers**

► 1−**norm** vs 2−**norm**

- as usual, the 1−**norm** tends to <u>**limit**</u> **more** drastically the **outlier contributions**

- this makes it a bit **more robust**, and it tends to be used <u>more</u> frequently in practice

► <u>common problem</u>:

- **not** really intuitive <u>**how to set up** $C$</u>

- usually **cross−validation**: there is a need to cross−validate with respect to both $C$ and kernel parameters

# $v - $ **SVM**

▶ an **alternative** formulation has been introduced to try to **overcome** this

$$\min_{\mathbf{w}, \boldsymbol{\xi}, \rho, b} \ \|\mathbf{w}\|^2 - v\rho + \frac{1}{n}\sum_i \xi_i$$

$$\text{subject to} \quad y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq \rho - \xi_i, \forall i$$

$$\xi_i \geq 0, \forall i$$

$$\rho \geq 0, \forall i$$

▶ **advantages:**

- $v$ has **intuitive** interpretation:

  1) $v$ is an <u>upper bound</u> on the proportion of training vectors that are margin errors, i.e. for which $y_i g(\mathbf{x}_i) \leq \rho$

  2) $v$ is a <u>lower bound</u> on total number of support vectors

- more discussion on Quiz #4 (Prob. 3)

# SVM: Connections to Regularization

▶ we talked about <u>penalizing</u> functions that are <u>too</u> complicated to improve **generalization**

▶ instead of the empirical risk, we should minimize the <u>regularized risk</u>

$$R_{reg}[f] = R_{emp}[f] + \lambda\Omega[f]$$

$$\min_{\mathbf{w},\xi,b} \; \|\mathbf{w}\|^2 + Cf(\xi)$$
$$\text{subject to} \quad y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$$
$$\xi_i \geq 0, \forall i$$

▶ the SVM seems to be <u>**doing this**</u> in some sense:

- it is designed to have as <u>few</u> errors as possible on **training set** (this is **controlled** by the **soft−margin weight** $C$)

- we <u>maximize</u> the margin by **minimizing $\|\mathbf{w}\|^2$** (which is a form of **complexity penalty**)

- hence, **maximizing the margin** must be connected to <u>enforcing</u> <u>some form of regularization</u>

# SVM: Connections to Regularization

- the connection can be made <u>explicit</u>

- consider the $1-$norm SVM

$$\min_{\mathbf{w},\boldsymbol{\xi},b} \|\mathbf{w}\|^2 + C \sum_i \xi_i \quad \text{subject to} \quad y_i g(\mathbf{x}_i) \geq 1 - \xi_i, \forall i$$
$$\xi_i \geq 0, \forall i$$

- the constraints can be rewritten as

$$\text{i) } \xi_i \geq 0 \quad \text{and} \quad \text{ii) } \xi_i \geq 1 - y_i g(\mathbf{x}_i)$$

which is equivalent to

$$\xi_i \geq \max[0, 1 - y_i g(\mathbf{x}_i)] = [1 - y_i g(\mathbf{x}_i)]_+$$

- note that the cost $\|\mathbf{w}\|^2 + C \sum_i \xi_i$ can only **increase** with larger $\xi_i$

- hence, at the optimal solution, $\quad \xi_i^* = [1 - y_i g(\mathbf{x}_i)]_+$

# SVM: Connections to Regularization

► the problem   $\xi_i^* = \max[0, 1 - y_i g(\mathbf{x}_i)] = [1 - y_i g(\mathbf{x}_i)]_+$

$$\min_{\mathbf{w},b} \|\mathbf{w}\|^2 + C \sum_i [1 - y_i g(\mathbf{x}_i)]_+ \iff \min_{\mathbf{w},b} \sum_i [1 - y_i g(\mathbf{x}_i)]_+ + \lambda \|\mathbf{w}\|^2$$

(by making $\lambda = 1/C$)

can be seen as a

## regularized risk

$$R_{reg}[f] = \sum_i L[\mathbf{x}_i, y_i, f] + \lambda \Omega[f]$$

with

- loss function   $L[\mathbf{x}, y, g] = [1 - y g(\mathbf{x})]_+$  →  hinge loss

- standard regularizer   $\Omega[\mathbf{w}] = \|\mathbf{w}\|^2$

# The SVM Loss

▶ it is interesting to **compare** the SVM loss

$$L[\mathbf{x}, y, g] = [1 - yg(\mathbf{x})]_+$$

with the 0/1 loss:

- the SVM loss **penalizes large** negative margins

- assigns some penalty to anything with margin less than 1

- for the 0/1 loss, the errors are **all** the same

SVM loss

0/1

1

0     1     $yg(x)$

▶ the **regularizer**

$$\Omega[\mathbf{w}] = \|\mathbf{w}\|^2$$

- **penalizes** planes of **large** $\|\mathbf{w}\|$

- standard **measure of complexity** in regularization theory

# Recap: Risk Minimization

▶ note that **all** the methods we have studied minimize a <u>similar risk</u>
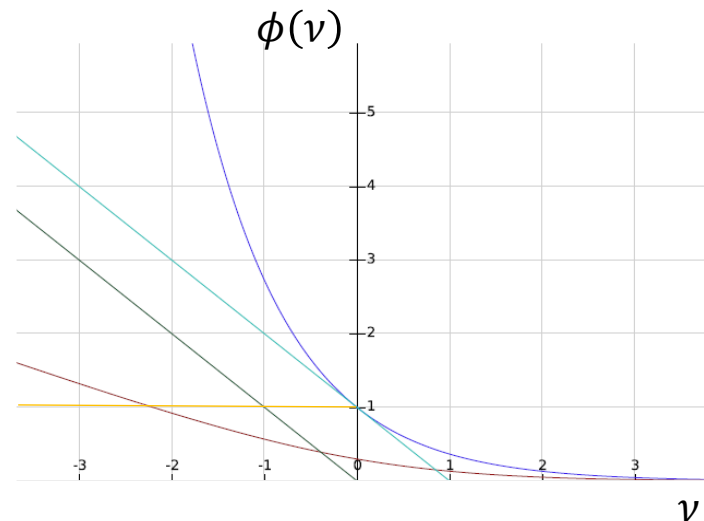
$$R_{reg}[f] = \sum_i L[y_i, f(\mathbf{x}_i)] + \lambda \Omega[f]$$

▶ in **all** cases, the <u>loss function</u> $L[y, g(\mathbf{x})]$ is a <u>margin loss</u>

$$L[y, g(\mathbf{x})] = \phi(yg(\mathbf{x}))$$

▶ only **difference** is the $\phi(\cdot)$ function

| Method | $\phi(v)$ | |
|---|---|---|
| BDR (0/1 loss) | $\text{sign}(-v)$ | 🟨 |
| Perceptron | $[-v]_+$ | 🟩 |
| neural networks | $\log(1 + e^{-v})$ | 🟫 |
| boosting | $e^{-v}$ | 🟦 |
| SVM | $[1-v]_+$ | 🔵 |

# Recap: Risk Minimization

▶ note that **all** the methods we have studied minimize a <u>similar risk</u>

$$R_{reg}[f] = \sum_i L[y_i, f(\mathbf{x}_i)] + \lambda\Omega[f]$$

▶ the **regularizer** $\Omega[f]$ is implemented in **different** ways

| Method | $\Omega[f]$ |
|---|---|
| BDR (0/1 loss) | enforced in the estimation of the pdfs |
| Perceptron | none |
| neural networks | weight decays |
| boosting | regularization is implemented by limiting the number of iterations (weak learners) |
| SVM | $\|\mathbf{w}\|^2$ |

# Recap: Risk Minimization

$$R_{reg}[f] = \sum_i L[y_i, f(\mathbf{x}_i)] + \lambda\Omega[f]$$

▶ note that minimizing $R_{reg}$ is the same as maximizing

$$e^{-R_{reg}[f]} = e^{-\sum_i L[y_i, f(\mathbf{x}_i)]} \cdot e^{-\lambda\Omega[f]}$$

which is the **same** as

- finding the function $f$ of <u>**maximum**</u> a posteriori probability
- under a probabilistic model with

likelihood function

$$e^{-\sum_i L[y_i, f(\mathbf{x}_i)]}$$

prior

$$e^{-\lambda\Omega[f]}$$

▶ hence, it has a **Bayesian interpretation**, where the **regularizer** defines the **prior**, which is used to constrain the values of the solution
(e.g., if $\Omega[\mathbf{w}] = \|\mathbf{w}\|^2$, the prior will be Gaussian with $\boldsymbol{\mu} = \mathbf{0}$ and $\boldsymbol{\Sigma} = \mathbf{I}$)

# In Summary

▶ <u>all</u> methods are implementations of the same optimization framework

▶ loss functions can have **<u>significant difference</u>** (margin enforcing vs not)

▶ regularizers are more tied to the <u>implementation</u>