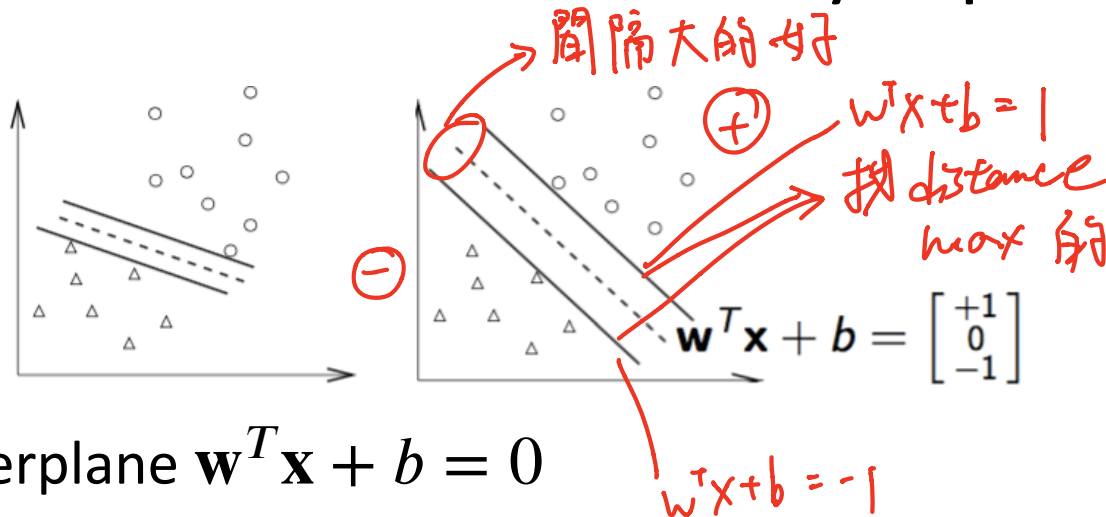# Support vector machines

## Hung-Hsuan Chen

Many are taken from Prof. C.-J. Lin's and J. Leskovec's slides

# (Linear) support vector classification

- Data point $i$: $\mathbf{x}_i = (x_{i1}, x_{i2}, \ldots, x_{id})$

- Class label of $i$: $y_i$

  - Two classes
  - Class 1: $y_i = 1$
  - Class 2: $y_i = -1$

- Find a hyperplane to separate the data points

# Assume the dataset is linearly separable

間隔大的好

$+$   $w^T x + b = 1$

找 distance max 的

$w^T x + b = \begin{bmatrix} +1 \\ 0 \\ -1 \end{bmatrix}$

$-$

$w^T x + b = -1$

- A hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$

  $\mathbf{w}^T \mathbf{x}_i + b \geq 1$ if $y_i = 1$

  $\mathbf{w}^T \mathbf{x}_i + b \leq -1$ if $y_i = -1$

- Discriminant function $f(\mathbf{x}) = \mathrm{sgn}(\mathbf{w}^T \mathbf{x} + b)$

  – There are **many different choices** of **w** and b

# Margin distance

- Given two parallel hyperplanes $H_1$ and $H_2$

$$H_1 : \mathbf{w}^T \mathbf{x} = b_1$$
$$H_2 : \mathbf{w}^T \mathbf{x} = b_2$$

- The distance between $H_1$ and $H_2$ is

$$d\left(H_1, H_2\right) = \frac{|b_1 - b_2|}{||\mathbf{w}||_2}$$

- Distance between $\mathbf{w}^T \mathbf{x}_i + b = 1$ and $\mathbf{w}^T \mathbf{x}_i + b = -1$:

$$\text{margin} = \frac{2}{||\mathbf{w}||_2}$$

# Maximum margin

- $\mathbf{w}, b = \mathrm{argmax}_{\mathbf{w},b} \dfrac{2}{\left|\left|\mathbf{w}\right|\right|_2}$

- This is the same as

$$\mathbf{w}, b = \mathrm{argmin}_{\mathbf{w},b} \frac{1}{2}\mathbf{w}^T\mathbf{w}$$

- This is modeled as a **quadratic programming** problem

$$\min_{\mathbf{w},b} \frac{1}{2}\mathbf{w}^T\mathbf{w}$$

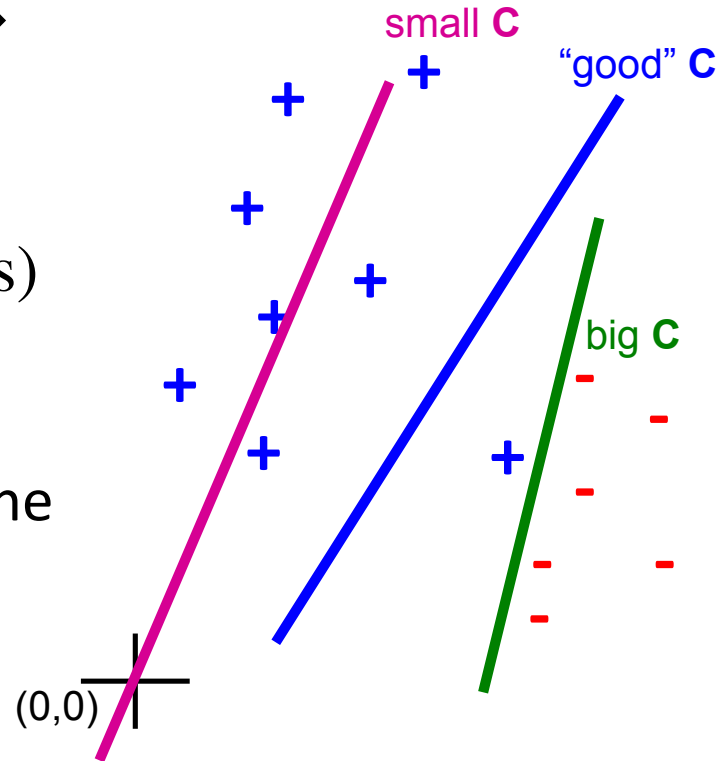Subject to $y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right) \geq 1 \ \forall i$

# Non-linearly separable dataset

- If non-linearly separable ➔ introduce penalty

$$\min_{\mathbf{w},b} \frac{1}{2}\mathbf{w}^T\mathbf{w} + C(\text{\# of mistakes})$$

  - If $C \to \infty$: allows no error
  - If $C = 0$: basically ignores the data at all

small **C**

"good" **C**

big **C**

+ + + + + + + + +

- - - - - -
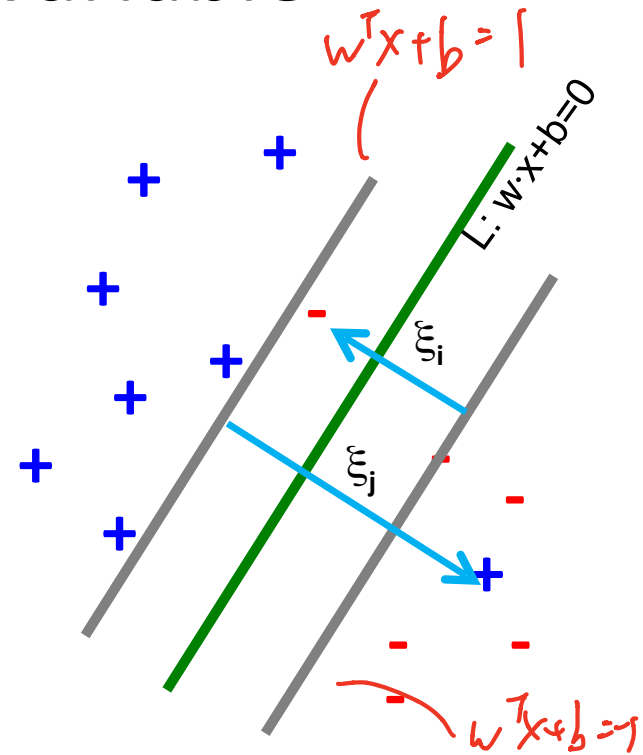
(0,0)

# Introduce slack variable

- Not all mistakes are equally bad

$$\min_{\mathbf{w},b} \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{n} \xi_i$$

Subject to

$$y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right) \geq 1 - \xi_i \ \forall i$$

- If a point is on the wrong side ➜ get penalty $\xi_i$

$w^Tx+b=1$

L: w·x+b=0

$\xi_i$

$\xi_j$

$w^Tx+b=-1$

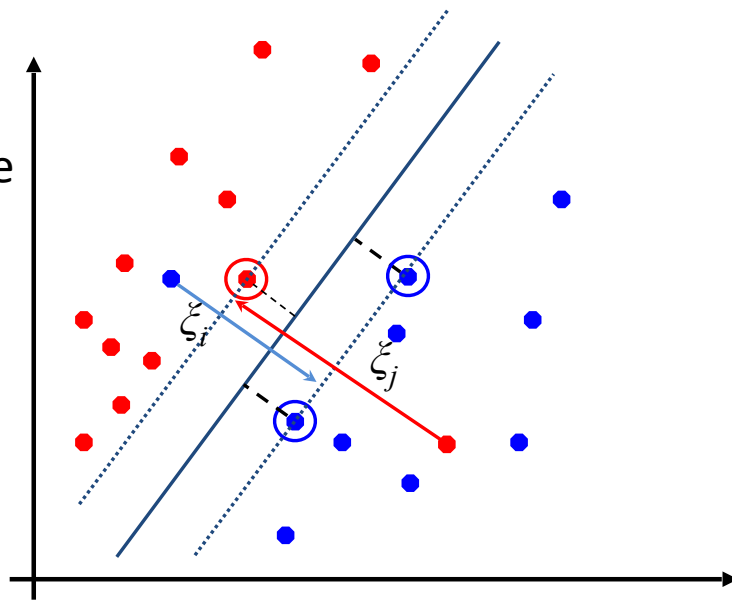**For each data point x:**
If d(x, L) ≥ 1 and at the right side: don't care
Else: pay linear penalty

# Soft margin classification

- ## Why soft margin
  - The training data may not be linearly separable
  - Even if the training data is linearly separable, allowing some error may increase the margin
- ## Essentially, there are two objectives (which may against each other)
  - Minimize the training error
    - Prevent error
  - Maximize the margin
    - Prevent overfitting (allow some error)

# Soft margin classification formula

- Original (linear) formula

$$\min_{\mathbf{w},b} \frac{1}{2}\mathbf{w}^T\mathbf{w}$$

Subject to

$$y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right) \geq 1 \;\; \forall i$$

- New formula

$$\min_{\mathbf{w},b}\left( \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_i \xi_i \right)$$

Subject to

$$y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right) \geq 1 - \xi_i$$

and $\xi_i > 0 \;\; \forall i$

- $C$: control overfitting
  - A large $C$ makes most $\xi_i$'s to zero
- $\xi_i$: slack variables

# Linear SVM with soft margin

- Linear SVM

$$\min_{\mathbf{w},b} \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{n}\xi_i$$

Subject to

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

- This is the same as

$$\min_{\mathbf{w},b}\left(\frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{n}\max\left\{0, 1 - y_i(\mathbf{w}^T\mathbf{x}_i + b)\right\}\right)$$

*handwritten annotations:*

⊗ if data incorrect predicted
1−0 would be positive
would be 1−0 (>0)

∅ if data correct predicted
⟶ 1−0 would be negative
would be 0

If the point is at the wrong side, get loss proportional to $\xi_i$

Margin inverse

Regularization parameter

Empirical **loss L** (how well we fit training data)

$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

縮後可用 GD

# Derivatives

$$f(\mathbf{w}, b) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{n}\max\left\{0,\, 1 - y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right)\right\}$$

$$\Rightarrow \nabla_{w_j}f = w_j + C\sum_{i=1}^{n}\frac{\partial\max\left\{0,\, 1 - y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right)\right\}}{\partial w_j} = \begin{cases} w_j & \text{if } y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right) \geq 1 \\ w_j + C\left(-y_i x_{ij}\right) & \text{else} \end{cases}$$

means correct predicted

cause: $1 - y_i(w^T x_i + b) < 0$

# Solve Linear SVM by GD

```
While (true) {
    for (j=1,2, ..., d){
```
*(d crossed out, replaced with N above)*

Note: $b$ is batch size *(crossed out)*

$$\nabla_{w_j} f\left(\mathbf{x}_{1:d}\right) = w_j + C \sum_{i=1}^{d} \frac{\partial \max\left\{0,\, 1 - y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right)\right\}}{\partial w_j}$$

*(subscript $d$ crossed out, replaced with $N$; summation upper limit $d$ crossed out, replaced with $N$)*

$$w_j = w_j - \alpha \nabla_{w_j} f$$

```
    }

    if (w converges) break
}
```

# Solve Linear SVM by SGD

```
for (i=1,2, ..., n){
  for (j=1,2, ..., d){
```

$$\nabla_{w_j} f(\mathbf{x}_i) = w_j + C \frac{\partial \max\left\{0,\, 1 - y_i\left(\mathbf{w}^T \mathbf{x}_i + b\right)\right\}}{\partial w_j}$$
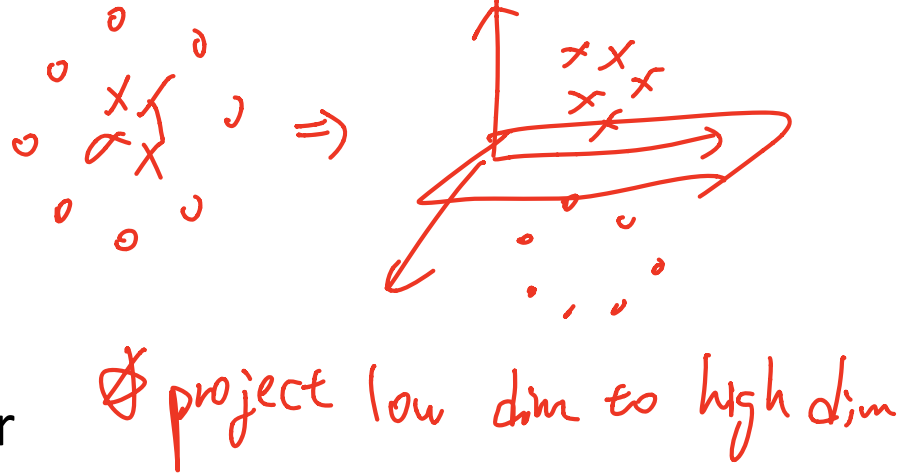
$$w_j = w_j - \alpha \nabla_{w_j} f$$

```
  }

  if (w converges) break
}
```

$2D \rightarrow 3D$

- Detour
  - Lagrange multiplier
  - KKT condition

- Math caution!
  - If you get lost, I hope you at least understand the linear SVM

*project low dim to high dim*

→ use to check Lagrange 是否符合 Solution!

# Generalized Lagrange multiplier

- Standard form problem

  **Minimize** $f(\mathbf{x})$ subject to $g_i(\mathbf{x}) \leq 0$ $(i = 1, \ldots, p)$

  and $h_j(\mathbf{x}) = 0$ $\left(j = 1, \ldots, m\right)$

- Lagrangian

$$\mathscr{L}\left(\mathbf{x}, \lambda, \mu\right) = f(\mathbf{x}) + \sum_{i=1}^{p} \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^{m} \mu_j h_j(\mathbf{x})$$

# The characteristic of the solution

SVM (ignore slack variables):
$$\min_{\mathbf{w},b} \frac{1}{2}\mathbf{w}^T\mathbf{w}$$
Subject to
$$y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 \ \forall i$$

$$\mathcal{L}(\mathbf{w}, b, \lambda) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + \sum \lambda_i\left[1 - y_i(\mathbf{w}^T\mathbf{x}_i + b)\right]$$

– No equality constraints (no μ's)

- Based on the KKT condition: if $\mathbf{w}^*$ is the optimal solution to the standard form problem, then there exist KKT multipliers $\lambda$ and $\mu$ such that

  – Lagrangian optimality
  $$\nabla\mathcal{L}(\mathbf{w}^*, \lambda, \mu) = 0 - - - - - \quad (1)$$

  – Primal feasibility
  $$g_i(\mathbf{w}^*) \leq 0 \ \forall i - - - - - - - \quad (2)$$
  $$h_j(\mathbf{w}^*) = 0 \ \forall j - - - - - - - \quad (3)$$

  – **Dual feasibility**
  $$\lambda_i \geq 0 \ \forall i - - - - - - - - - \quad (4)$$

  – **Complementary slackness**
  $$\lambda_i g_i(\mathbf{w}^*) = 0 \ \forall i - - - - - \quad (5)$$

*(handwritten annotations)* $w^Tx + b = -1$ ; $w^Tx + b = 1$ ; $w^Tx+b \geq 1$ ; $w^Tx+b \leq -1$ ; $\ominus$ ; $\oplus$

$$\mathbf{w}^T\mathbf{x}_i + b = 1 \text{ and } \mathbf{w}^T\mathbf{x}_i + b = -1$$

*(handwritten)* 只有线上的 $g(w) = 0$
以外的 $\lambda = 0$，所以 $g() \neq 0$

- Assume linearly-separable, by condition (4) and (5):

  – If a training instance $\mathbf{x}_i$ is **_not_** on the two hyperplanes (i.e., $g_i(\mathbf{w}^*) < 0$), $\lambda_i$ must be 0

  – If a training instance $\mathbf{x}_i$ is on the two hyperplanes (i.e., $g_i(\mathbf{w}^*) = 0$), $\lambda_i \geq 0$

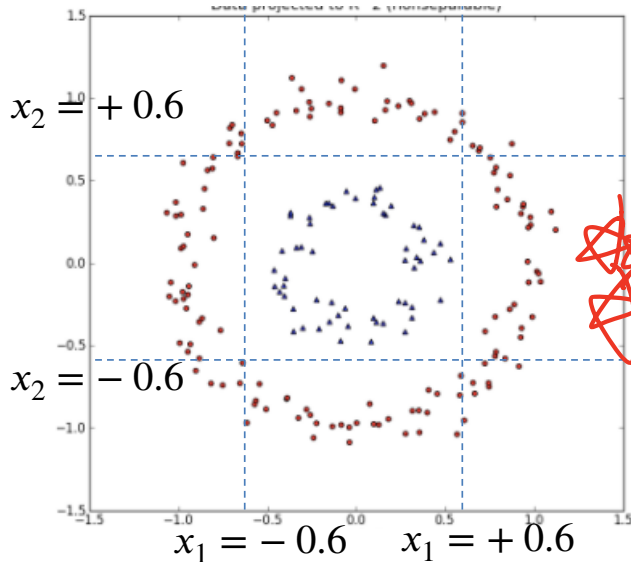# Map features to higher dimensional



kernel

Decision surface

- Transform the data into a higher dimension feature space so that linear separation is possible
  - Higher dimensional (**could be infinite**) feature space

$$\bullet \ \phi(\mathbf{x}_i) = \left[\phi_1(\mathbf{x}_i), \phi_2(\mathbf{x}_i), \ldots\right]^T$$

2D to 3D

# Example

$x_2 = +0.6$

$x_2 = -0.6$

$x_1 = -0.6$    $x_1 = +0.6$

- The positive and negative examples are not linearly- separable by the 2D features $(x_1, \ x_2)$
- If we add one more feature $x_3 = x_1^2 + x_2^2$, the blue points are those with $x_3 \leq 0.6^2$, and the red points are those with $x_3 > 0.6^2$
  - $\phi(x_1, x_2) = (x_1, x_2, x_3) = (x_1, x_2, x_1^2 + x_2$
    - 2D to 3D
  - The points become linearly separable

11/3/20
18

# Kernel SVM

$$\min_{\mathbf{w},b} \frac{1}{2}\mathbf{w}^T\mathbf{w}$$

Subject to $y_i\left(\mathbf{w}^T \phi(\mathbf{x}_i) + b\right) \geq 1$

Here we ignore the slack variables for simplicity

- Linear SVM: length of **w** is $d$ (the same as the size of $\mathbf{x}_i$)
- Kernel SVM: length of **w** is larger than $d$ (the same as the size of $\phi\left(\mathbf{x}_i\right)$)
  - Kernel SVM can fit a more complex function
    - The size of $\phi\left(\mathbf{x}_i\right)$ is large (and could be *infinity*)
  - How to efficiently compute **w** and $\mathbf{w}^T\phi(\mathbf{x}_i)$?
  - How to store **w**?

# Lagrangian of kernel SVM

$g(\ )$

$$\mathscr{L}\left(\mathbf{w}, b, \boldsymbol{\lambda}\right) = \frac{1}{2}\mathbf{w}^{T}\mathbf{w} + \sum \lambda_i \left[1 - y_i\left(\mathbf{w}^{T}\phi(\mathbf{x}_i) + b\right)\right]$$

$$\begin{cases} \nabla_{\mathbf{w}}\mathscr{L}\left(\mathbf{w}, b, \boldsymbol{\lambda}\right) = \mathbf{w} - \sum \lambda_i y_i \phi(\mathbf{x}_i) := 0 \\ \nabla_{b}\mathscr{L}\left(\mathbf{w}, b, \boldsymbol{\lambda}\right) = -\sum \lambda_i y_i := 0 \end{cases}$$

$$\Rightarrow \begin{cases} \mathbf{w} = \sum \lambda_i y_i \phi(\mathbf{x}_i) \\ \sum \lambda_i y_i := 0 \end{cases}$$

- Given a test instance $\mathbf{x}_t$, the discriminant function is

$$f\left(\mathbf{x}_t\right) = \mathbf{w}^{T}\phi(\mathbf{x}_t) + b = \sum \lambda_i y_i \phi(\mathbf{x}_i)^{T}\phi(\mathbf{x}_t) + b$$

  - Prediction is a **linear combination of training instances** $\mathbf{w}^{T}\phi(\mathbf{x}_t)$ plus bias

# High dimensional mapping example

$$f(\mathbf{x}_t) = \mathbf{w}^T \phi(\mathbf{x}_t) + b = \sum \lambda_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_t) + b$$

- Example:
  - $\mathbf{x}_i = [x_{i1}, x_{i2}]^T \in R^2, \phi(\mathbf{x}_i) \in R^6$
  - If we set (assume)
  
  $$\phi(\mathbf{x}_i) = [1, \sqrt{2}x_{i1}, \sqrt{2}x_{i2}, \sqrt{2}x_{i1}x_{i2}, \ x_{i1}^2, x_{i2}^2]^T$$
  
  - Then
  
  $$\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_t) = 1 + x_{i1}^2 x_{t1}^2 + x_{i2}^2 x_{t2}^2 + 2x_{i1}x_{t1} + 2x_{i2}x_{t2} + 2x_{i1}x_{t1}x_{i2}x$$
  
  – When the target dimension is large, it is inefficient to generate $\phi(\mathbf{x}_t)$ and $\phi(\mathbf{x}_i) \ \forall i$ and perform the dot product

# Kernel trick example

- If

$$\phi(\mathbf{x}_i) = [1, \sqrt{2}x_{i1}, \sqrt{2}x_{i2}, \sqrt{2}x_{i1}x_{i2},\ x_{i1}^2, x_{i2}^2]^T,$$

then:

  - $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_t) = \left(1 + \mathbf{x}_i^T \mathbf{x}_T\right)^2$

  - Computing $\left(1 + \mathbf{x}_i^T \mathbf{x}_T\right)^2$ is much more efficient than computing $\phi(\mathbf{x}_i)$, $\phi(\mathbf{x}_t)$, and then $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_t)$

# Kernel trick example

$$f(\mathbf{x}_t) = \mathbf{w}^T \phi(\mathbf{x}_t) + b = \sum \lambda_i y_i \boxed{\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_t)} + b$$

所以計算π

- If $\phi(\mathbf{x}_i)$'s dimension is very high
  - Store $\mathbf{w}$ is costly
  - Compute discriminant function $f(\mathbf{x}_t) = \mathbf{w}^T \phi(\mathbf{x}_t) + b$ is costly

- We may use $\left(1 + \mathbf{x}_i^T \mathbf{x}_T\right)^2$ to efficiently map features to higher dimension

- We compute

$$\sum \lambda_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_t) + b = \sum \lambda_i y_i \left(1 + \mathbf{x}_i^T \mathbf{x}_T\right)^2 + b \text{ as the}$$

discriminant function

☆ most of data are 0
只有在線上的才 > 0

# Popular kernels

- Linear kernel (i.e., linear SVM)

$$K\left(\mathbf{x}_i, \mathbf{x}_t\right) = \mathbf{x}_i^T \mathbf{x}_t = \left\langle \mathbf{x}_i, \mathbf{x}_t \right\rangle$$

- Polynomial kernel

$$K\left(\mathbf{x}_i, \mathbf{x}_t\right) = \left(\left\langle \mathbf{x}_i, \mathbf{x}_t \right\rangle + r\right)^d, \ r > 0$$

- Gaussian (RBF) kernel

$$K\left(\mathbf{x}_i, \mathbf{x}_t\right) = \exp\left(-\gamma \|\mathbf{x}_i - \mathbf{x}_t\|^2\right)$$

- The dimension of $K\left(\mathbf{x}_i, \mathbf{x}_t\right)$ could be ***infinity*** (e.g., RBF kernel), but the dimensions of $\mathbf{x}_i$ and $\mathbf{x}_t$ are finite

# Mapping to infinite dimensional

- Assume $\mathbf{x}_i \in R^1$, $\gamma > 0$
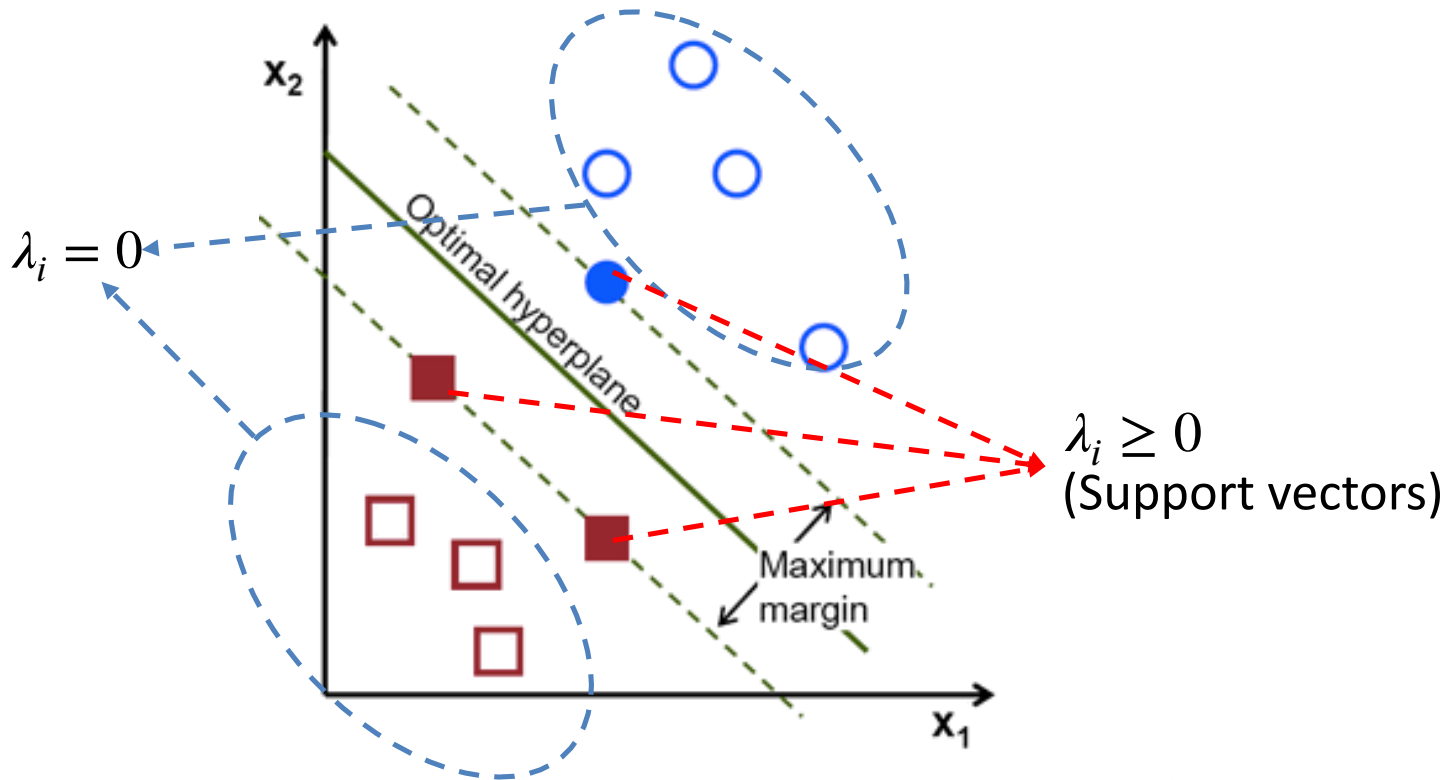
$$\exp\left(-\gamma\|\mathbf{x}_i - \mathbf{x}_t\|^2\right) = \exp\left(-\gamma(\mathbf{x}_i - \mathbf{x}_t)^2\right) = \exp\left(-\gamma\mathbf{x}\right.$$

$$= \exp\left(-\gamma\mathbf{x}_i^2 - \gamma\mathbf{x}_j^2\right) \cdot \exp\left(2\gamma\,\mathbf{x}_i\mathbf{x}_j\right)$$

$$= \exp\left(-\gamma\mathbf{x}_i^2 - \gamma\mathbf{x}_j^2\right)\left[1 + \frac{2\gamma\,\mathbf{x}_i\mathbf{x}_j}{1!} + \frac{\left(2\gamma\,\mathbf{x}_i\mathbf{x}_j\right)^2}{2!} + \frac{\left(2\gamma\,\mathbf{x}_i\mathbf{x}_j\right)^3}{3!} + \cdots\right]$$

$$=$$

$$\exp\left(-\gamma\mathbf{x}_i^2 - \gamma\mathbf{x}_j^2\right)\left(1 \cdot 1 + \sqrt{\frac{2\gamma}{1!}}\mathbf{x}_i\sqrt{\frac{2\gamma}{1!}}\mathbf{x}_j + \sqrt{\frac{(2\gamma)^2}{2!}}\mathbf{x}_i^2\sqrt{\frac{(2\gamma)^2}{2!}}\mathbf{x}_j^2 + \sqrt{\frac{(2\gamma)^3}{3!}}\mathbf{x}_i^3\sqrt{\frac{(2\gamma)^3}{3!}}\mathbf{x}_j^3 + \cdots\right)$$

$$= \phi\left(\mathbf{x}_i\right)^T \phi\left(\mathbf{x}_j\right),$$

Where $\phi\left(\mathbf{x}_i\right) = \exp\left(-\gamma\mathbf{x}_i^2\right)\left[1, \sqrt{\frac{2\gamma}{1!}}\mathbf{x}_i, \sqrt{\frac{(2\gamma)^2}{2!}}\mathbf{x}_i^2, \sqrt{\frac{(2\gamma)^3}{3!}}\mathbf{x}_i^3, \cdots\right]^T$

# Characteristics of the solution

$$\to \Sigma \lambda_i y_i \phi(x)^{\mathsf{T}}$$

- Discriminant function

$$f(\mathbf{x}_t) = \mathbf{w}^T \phi(\mathbf{x}_t) + b = \sum \lambda_i y_i K(\mathbf{x}_i, \mathbf{x}_t) + b$$

$$\phi(x_i)^{\mathsf{T}} \phi(x_t)$$

- Many $\lambda_i$'s are 0
  - Memorizing training instance $(\mathbf{x}_i, y_i)$ only if $\lambda_i > 0$
  - We don't need to form **w** explicitly

- To predict the label of a test instance $\mathbf{x}_t$, we need to compute the ***Kernel*** of the test instance with the training instances **whose $\lambda_i$'s are larger than zeros**
  - These training instances are called "**support vectors**"

# Visualizing "support vectors"



$\lambda_i = 0$

$\lambda_i \geq 0$
(Support vectors)

Optimal hyperplane

Maximum margin

$x_2$

$x_1$

11/3

# A numerical example

- Training data: five 1D data points with labels
  - First class (+1): $x_{1,1} = 1$, $x_{2,1} = 2$, $x_{5,1} = 6$
  - Second class (-1): $x_{3,1} = 4$, $x_{4,1} = 5$
- Non-linearly separable
- Use polynomial kernel with degree 2

$$K(\mathbf{x}_i, \mathbf{x}_t) = \left( \langle \mathbf{x}_i, \mathbf{x}_t \rangle + 1 \right)^2$$
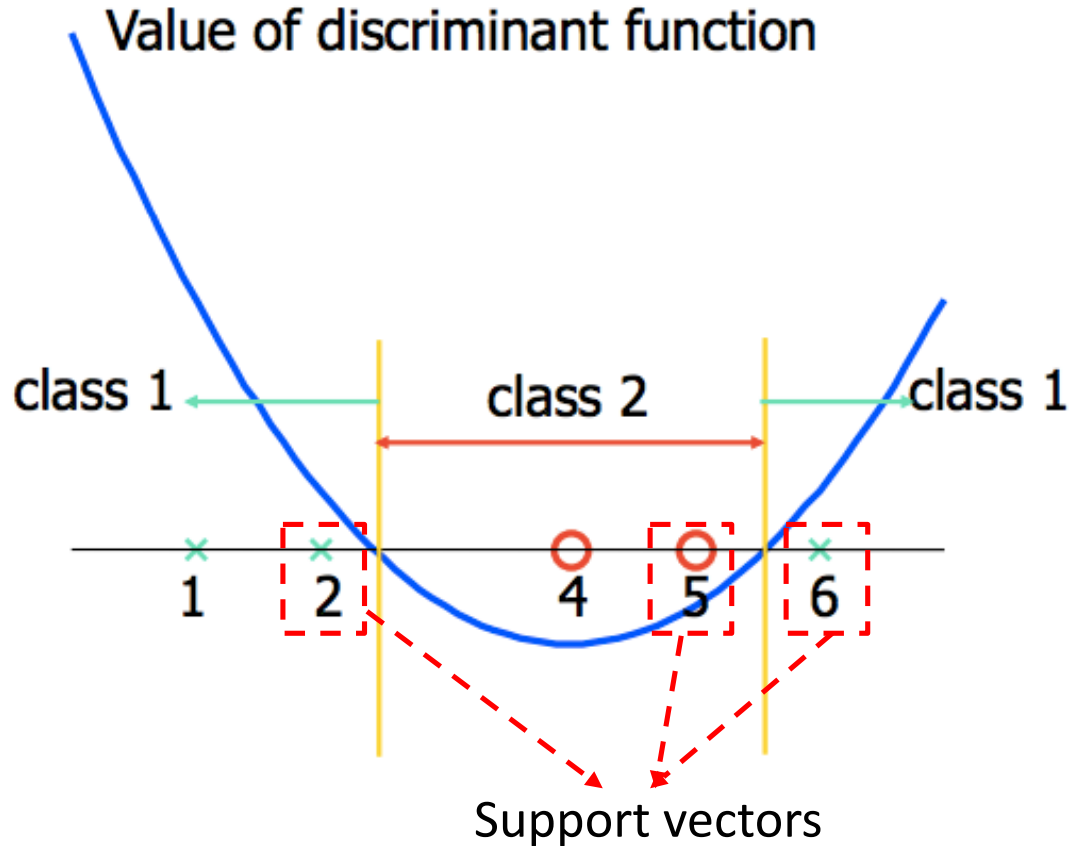
- Solve $\lambda_i$ $(i = 1, \ldots, 5)$ by a standard QP solver
  - $\left[ \lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5 \right] = [0,\ 2.5,\ 0,\ 7.333,\ 4.833]$

# The discriminant function of the example

- $\left[\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5\right] = [0,\, 2.5,\; 0,\, 7.333,\; 4.833]$
  - Since $\lambda_1 = \lambda_3 = 0$, the support vectors are
  $$\left(x_2,\; x_4,\; x_5\right) = (2,\, 5,\, 6)$$
- The discriminant function

$$\circ \; f(z) = \sum \lambda_i y_i \phi\left(x_i\right)^T \phi(z) + b$$
$$= \left[ 2.5(1)(2z+1)^2 + 7.333(-1)\left(5z+1\right)^2 \right.$$
$$\left. +4.833(1)\left(6z+1\right)^2 \right] + b$$
$$= 0.6667z^2 - 5.333z + b$$

$\circ$ Since $f\left(x_{2,1}\right) = f\left(x_{5,1}\right) = 1$ and $f\left(x_{4,1}\right) = -1$, we can get $b$=9

$\circ\; f(z) = 0.6667z^2 - 5.333z + 9$

# Visualize the discriminant function

# Standard SVM

- Standard form

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

Subject to (1) $y_i\left(\mathbf{w}^T\phi(\mathbf{x}_i) + b\right) \geq 1 - \xi$

(2) $\xi_i > 0 \;\; \forall i$

# Review of SVM

- ## Large margin
  - Prevent overfitting

- ## Soft margin
  - Make the margin become larger
  - Prevent overfitting

- ## Kernel trick
  - Make the data linearly separable
  - Efficiently transform the input features into high (could be infinite) dimensional space

# Revisiting logistic regression and SVM from another viewpoint

# Deriving SVM and LogReg from regularized linear classification

- We derived SVM from the viewpoint of maximal margin

- We derived logistic regression from maximizing the log-likelihood (or minimizing the cross entropy loss)

- However, both can be considered from the viewpoint of ***regularized linear classification***

# Regularized linear classification

- Training data:

$$\{\mathbf{x}_i, y_i\}_{i=1,\ldots,n}, \ \mathbf{x_i} \in R^d, y_i \in \{\pm 1\}$$

- Objective

$$\min_{\mathbf{w}} f(\mathbf{w}), \ f(\mathbf{w}) \equiv \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^{n} \xi\left(\mathbf{w}; \mathbf{x}_i, y_i\right)$$

- $\xi\left(\mathbf{w}; \mathbf{x}_i, y_i\right)$: loss function; we hope $y_i \mathbf{w}^T \mathbf{x} > 0$
    - Trying to fit the training data
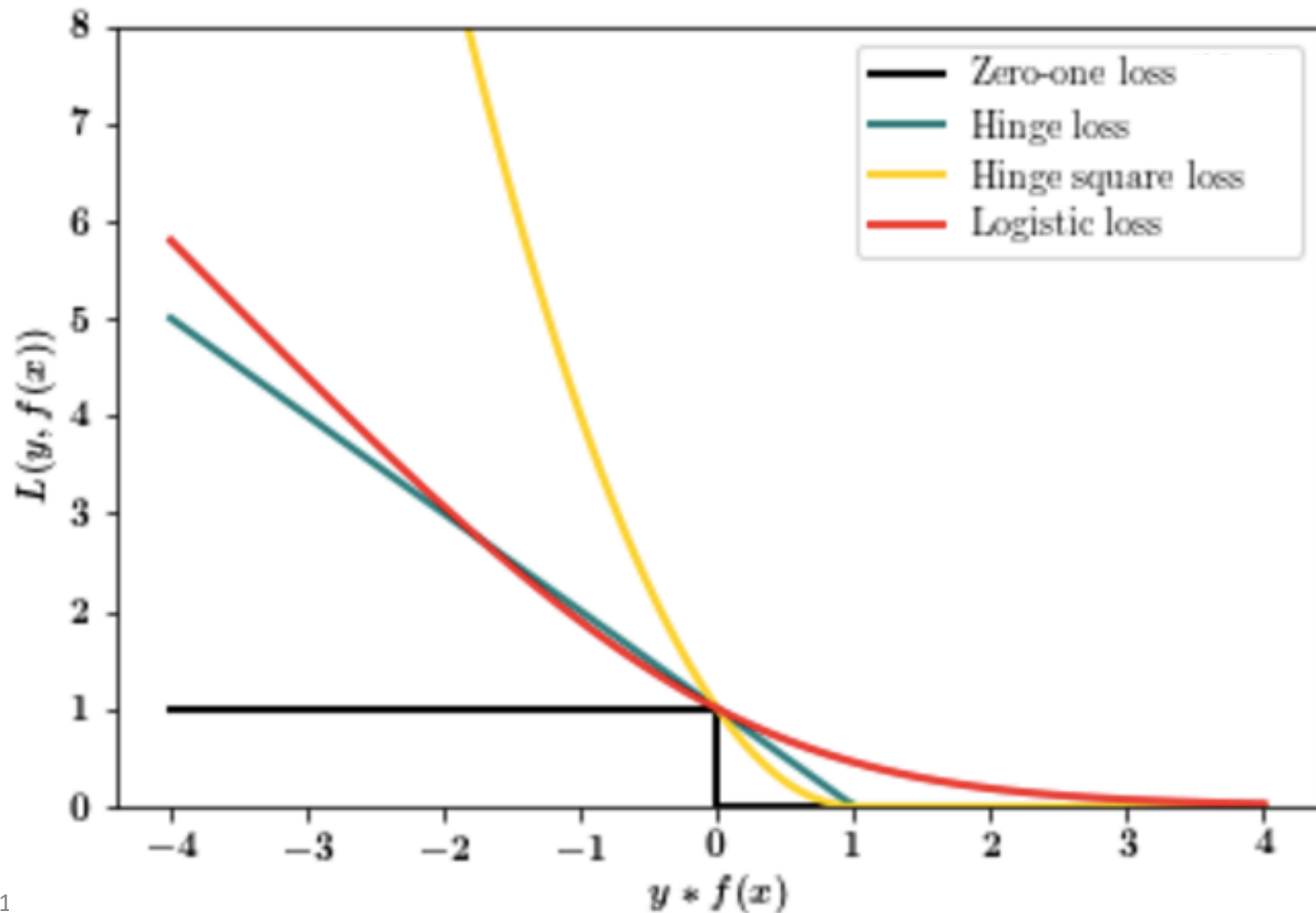- $\left(\mathbf{w}^T \mathbf{w}\right)/2$: regularization term
    - We skip the L1 regularization term here
    - Prevent over-fit the training data
- $C$: regularization parameter

# Loss functions

- Common loss functions in classification
  - Hinge loss
    - $\xi_{L1}\left(\mathbf{w}; \mathbf{x}_i, y_i\right) = \max\left(0,\ 1 - y_i \mathbf{w}^T \mathbf{x}_i\right)$
  - Squared hinge loss
    - $\xi_{L2}\left(\mathbf{w}; \mathbf{x}_i, y_i\right) = \max\left(0,\ 1 - y_i \mathbf{w}^T \mathbf{x}_i\right)^2$
  - Logistic loss
    - $\xi_{LR}\left(\mathbf{w}; \mathbf{x}_i, y_i\right) = \log\left(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}\right)$
      - This is different from what we derived previously
        » We used 1/0 to encode two classes before, but here we use +1/-1

- SVM: $\xi_{L1}$ and $\xi_{L2}$
- Logistic Regression: $\xi_{LR}$

# Visualizing the loss functions

# Regularized linear classification with kernel

- Training data:

$$\{\mathbf{x}_i, y_i\}_{i=1,\ldots,m}, \ \mathbf{x_i} \in R^n, y_i \in \{\pm 1\}$$

- Objective

$$\min_{\mathbf{w}} f(\mathbf{w}), f(\mathbf{w}) \equiv \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^{m} \xi\Big(\mathbf{w}; \phi(\mathbf{x}_i), y_i\Big)$$

  - $\xi\Big(\mathbf{w}; \phi(\mathbf{x}_i), y_i\Big)$: loss function; we hope $y_i \mathbf{w}^T \phi(\mathbf{x}_i) > 0$
    - Trying to fit the training data
  - $\mathbf{w}^T \mathbf{w}/2$: regularization term
    - We skip the L1 regularization term here
    - Prevent over-fit the training data
  - $C$: regularization parameter

# Logistic regression vs SVM

- Logistic regression and SVM are very related

- Their performance (i.e., test accuracy) is usually similar

- Due to the naming, the typical deriving process, and historical reasons, many believe that SVM and logistic regression are very different

  – This is a misunderstanding

# Linear or kernel? SVM or Logistic Regression?

- When people say SVM, they typically mean "**kernel** SVM"
  - But there is linear SVM

$$\min_{\mathbf{w}} f(\mathbf{w}), \ f(\mathbf{w}) \equiv \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^{m} \max\left(0, \ 1 - y_i \mathbf{w}^T \mathbf{x}_i\right)$$

- When people say logistic regression, they typically mean "**linear** logistic regression"
  - But there is kernel logistic regression

$$\min_{\mathbf{w}} f(\mathbf{w}), \ f(\mathbf{w}) \equiv \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^{m} \log\left(1 + e^{-y_i \mathbf{w}^T \phi(\mathbf{x}_i)}\right)$$

  - However, kernel logistic regression is rarely used in practice
    - Most $\lambda_i$ are not zero ➔ if we want to apply the kernel trick (instead of storing **w** explicitly), almost all training samples need to be memorized

# Regularized linear <span style="color:red">regression</span>

- Training data:

$$\left\{\mathbf{x}_i, y_i\right\}_{i=1,\ldots,m}, \ \mathbf{x_i} \in R^n, y_i \in R^1$$

- Objective

$$\min_{\mathbf{w}}\left( \frac{\mathbf{w}^T\mathbf{w}}{2} + C\sum_{i=1}^{m} \xi\left(\mathbf{w}; \mathbf{x}_i, y_i\right) \right)$$

  - $\xi\left(\mathbf{w}; \mathbf{x}_i, y_i\right)$: loss function
    - Trying to fit the training data
  - $\mathbf{w}^T\mathbf{w}/2$: regularization term
    - We skip the L1 regularization term here
    - Prevent over-fit the training data
  - $C$: regularization parameter

# Loss functions for regression

- Some commonly used loss functions
  - L1 loss
    - $\xi_{L1}\left(\mathbf{w}; \mathbf{x}_i, y_i\right) = \left| y_i - \mathbf{w}^T\mathbf{x}_i \right|$
  - L2 loss
    - $\xi_{L2}\left(\mathbf{w}; \mathbf{x}_i, y_i\right) = \left( y_i - \mathbf{w}^T\mathbf{x}_i \right)^2$
  - $\epsilon$-insensitive loss
    - $\xi_{\epsilon}\left(\mathbf{w}; \mathbf{x}_i, y_i\right) = \max\left( \left| \mathbf{w}^T\mathbf{x}_i - y_i \right| - \epsilon,\ 0 \right)$
  - $\epsilon$-insensitive square loss
    - $\xi_{\epsilon 2}\left(\mathbf{w}; \mathbf{x}_i, y_i\right) = \max\left( \left| \mathbf{w}^T\mathbf{x}_i - y_i \right| - \epsilon,\ 0 \right)^2$
- SVM (support vector regression): $\xi_{\epsilon}, \xi_{\epsilon 2}$
- Linear Regression: $\xi_{L2}$

# Regularized regression with kernel

- Training data:

$$\{\mathbf{x}_i, y_i\}_{i=1,\ldots,m}, \ \mathbf{x_i} \in R^n, y_i \in R^1$$

- Objective

$$\min_{\mathbf{w}} f(\mathbf{w}), f(\mathbf{w}) \equiv \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^{m} \xi\left(\mathbf{w}; \phi(\mathbf{x}_i), y_i\right)$$

  - $\xi\left(\mathbf{w}; \phi(\mathbf{x}_i), y_i\right)$: loss function
    - Trying to fit the training data
  - $\mathbf{w}^T\mathbf{w}/2$: regularization term
    - We skip the L1 regularization term here
    - Prevent over-fit the training data
  - $C$: regularization parameter

# Summary

- The same classification method can be derived from different ways
  - SVM
    - Maximize margin
    - Minimizing training loss with regularization constraints
  - LR
    - Maximize log-likelihood
    - Minimizing training loss with regularization constraints
- Linear regression and support vector regression are also under the same umbrella
- Understanding the concept of training loss and regularization enables you to self-study many machine learning techniques

# Quiz

- What are "support vectors" of SVM?

- When increasing training samples, will the size of "logistic regression model" increase?

- When increasing training samples, will the size of "linear SVM model" increase?

- When increasing training samples, will the size of "kernel SVM model" increase?

# Appendix

# Convex optimization and quadratic programming

- A convex optimization problem is one of the form

  Minimize $f_0(\mathbf{x})$

  Subject to $f_i(\mathbf{x}) \leq 0,\ i = 1, \ldots,\ m$

  $$\mathbf{a}_j^T \mathbf{x} = \mathbf{b}_j,\ j = 1, \ldots,\ p$$

  Where $f_0,\ \ldots,\ f_m$ are convex functions

- The convex optimization problem is called a quadratic programming (QP) if the objective function is (convex) quadratic, and the constraint functions are affine

  – Minimize $f_0(\mathbf{x}) = (1/2)\mathbf{x}^T\mathbf{P}\mathbf{x} + \mathbf{q}^T\mathbf{x} + \mathbf{r}$

  – Subject to $\mathbf{G}\mathbf{x} \preccurlyeq \mathbf{h}$ and $\mathbf{A}\mathbf{x} = \mathbf{b}$

# Primal problem

- The standard form problem re-formulate as the **Primal problem**

$$\min_{\mathbf{x}} \ \max_{\lambda_i \geq 0, \ \mu} \ \mathcal{L}(\mathbf{x}, \lambda, \mu)$$

- Why?

$$\max_{\lambda_i \geq 0, \ \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) = \max_{\lambda_i \geq 0, \ \mu} \left[ f(\mathbf{x}) + \sum_{i=1}^{p} \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^{m} \mu_j h_j(\mathbf{x}) \right]$$

$$= \max_{\lambda_i \geq 0, \ \mu} \left[ f(\mathbf{x}) + \sum_{i=1}^{p} \lambda_i g_i(\mathbf{x}) \right] \ \left( \because h_j(\mathbf{x}) = 0 \right)$$

$$= f(\mathbf{x}) \ \left( \because g_i(\mathbf{x}) \leq 0 \right)$$

$$\Rightarrow \min_{\mathbf{x}} \ \max_{\lambda_i \geq 0, \ \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) = \min_{\mathbf{x}} f(\mathbf{x})$$

# Dual problem

- Primal problem: $p^* = \min_{\mathbf{x}} \max_{\lambda_i \geq 0,\ \mu} \mathscr{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$

- Dual problem: $d^* = \max_{\lambda_i \geq 0,\ \mu} \min_{\mathbf{x}} \mathscr{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$

- $p^* \geq d^*$
  - The min of the max is no less than the max of the min
  - Duality gap: $p^* - d^*$

# Strong duality

- $d^* = p^*$
- If the following conditions are true, then strong duality holds
    1. $f$ and $g_i$'s are convex
    2. $h_i$'s are linear functions (i.e., Exists $a_i$ and $b_i$ such that $h_i(\mathbf{x}) = a_i^T \mathbf{x} + b_i$)
    3. Exists some $\mathbf{x}$ such that $g_i(\mathbf{x}) \leq 0$
- In SVM, the above conditions holds
    - We may solve the dual problem instead of the primal problem

# Lagrangian in SVM

$$\mathscr{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + \sum \lambda_i \left[ 1 - y_i(\mathbf{w}^T\mathbf{x}_i + b) \right]$$

— No equality constraints (no µ's)

$$\begin{cases} \nabla_{\mathbf{w}}\mathscr{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = \mathbf{w} - \sum \lambda_i y_i \mathbf{x}_i := 0 \\ \nabla_b \mathscr{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = -\sum \lambda_i y_i := 0 \end{cases}$$

$$\Rightarrow \begin{cases} \mathbf{w} = \sum \lambda_i y_i \mathbf{x}_i \\ \sum \lambda_i y_i := 0 \end{cases}$$

# Primal and dual problem in SVM

- Primal:

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

$$\begin{cases} \mathbf{w} = \sum \lambda_i y_i \mathbf{x}_i \\ \sum \lambda_i y_i := 0 \end{cases}$$

- Dual:

$$d^* = \max_{\lambda_i \geq 0} \min_{\mathbf{w}} \mathscr{L}(\mathbf{w}, \boldsymbol{\lambda}) = \max_{\lambda_i \geq 0} \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum \lambda_i \left[ 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) \right] = \max_{\lambda_i \geq 0} \min_{\mathbf{w}} \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum \left[ \lambda_i - \lambda_i y_i (\mathbf{w}^T \mathbf{x}_i + b) \right] = \max_{\lambda_i \geq 0} \left[ -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i \mathbf{x}_j \rangle + \sum \lambda_i \right]$$

Subject to: $\lambda_i \geq 0$ and $\sum \lambda_i y_i = 0$

∵ KKT condition    ∵ differentiate of the primal w.r.t. *b*

# The dual problem

$$\max_{\lambda_i \geq 0} \left[ -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i \mathbf{x}_j \rangle + \sum \lambda_i \right]$$

Subject to $\lambda_i \geq 0$ and $\sum \lambda_i y_i = 0$

- $\lambda_i$'s are the only unknowns
- This is a quadratic programming (QP) problem
  - A global maximum of $\lambda_i$'s can always be found
  - How to solve? Use standard QP solvers
- We don't compute explicitly compute $\mathbf{w}$
- Discriminant function:

$$f\left( \mathbf{x}_t \right) = \mathbf{w}^T \mathbf{x}_t + b = \sum \lambda_i y_i \langle \mathbf{x}_i, \mathbf{x}_t \rangle + b$$

# Optimization type of solving SVM

- Dual QP
  - SMO, SVM-light, etc.
  - Many available QP solvers
    - List: http://www.numerical.rl.ac.uk/people/nimg/qp/qp.html
- Primal SGD
  - NORMA
  - SVM-SGD
- Dual Coordinate Descent
  - LibLinear
- We skip the details here.  If interested, read
  - "Convex optimization" by Boyd and Vandenberghe
  - "Large-Scale Support Vector Machines: Algorithms and Theory" by Aditya Krishna Menon