CS2109S Tutorial 6 SVMs and Regularisation

(AY 24/25 Semester 2)

March 21, 2025

(Prepared by Benson)

# Contents

#### Regularisation

Recap Q1. L1-norm vs. L2-norm.

#### Support Vector Machines

Recap

- Q2. SVM
- Q3. Primal Formulation

#### **Bias and Variance**

Recap

Q4. Bias and Variance

#### Kernel Trick

Recap Bonus. Gaussian Kernel

# Recap: Regularisation

- Why do we need regularisation?
  - A. To decrease model complexity.
  - B. To improve training time efficiency.
  - C. To avoid underfitting.
  - D. To avoid overfitting.

# Recap: Regularisation

How does regularisation address overfitting?

- A. By penalizing weights for transformed features.
- B. By reducing noise in the training data.
- C. By penalizing large weights.
- D. By reducing the number of transformed features.

L2-norm (Ridge Regression):

$$J(\boldsymbol{w}) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\boldsymbol{w}}(\boldsymbol{x}^{(i)}) - y^{(i)})^2 \right] + \lambda \sum_{i=1}^{n} \boldsymbol{w}_i^2$$

L1-norm (Lasso Regression):

$$J(\boldsymbol{w}) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\boldsymbol{w}}(\boldsymbol{x}^{(i)}) - y^{(i)})^2 \right] + \lambda \sum_{i=1}^{n} |\boldsymbol{w}_i|^2$$



(a) No regularisation.



$$w_0=0.9, w_1=0.5$$
Cost  $pprox 0$ 

(b) L1 regularisation with  $\lambda = 5$ .



Total cost of the three points (from left to right):

- ▶ 4.2 + 1 = 5.2
- 2.2 + 2.5 = 4.7 <sup>1</sup>/<sub>2</sub>

▶ 
$$1.1 + 4 = 5.1$$

 $w_0 = 0.0, w_1 = 0.5$ Cost = 4.7

(c) L2 regularisation with  $\lambda = 5$ .



Total cost of the two points (from left to right):

▶ 0.5 + 2.6 = 3.1 <sup>™</sup>

$$w_0 = 0.2, w_1 = 0.25$$
  
Cost = 3.1



- L2 heavily penalizes larger parameters, preferring all smaller values.
- L1 may set values of certain parameters to 0 (why?).



If  $w_1$  is more important than  $w_0$ , pushing towards  $w_1$  is free lunch ("one for one").



Animation (See HTML slides):



 $\ell^1$  induces sparse solutions for least squares

L2-norm (Ridge Regression):

$$J(\boldsymbol{w}) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\boldsymbol{w}}(\boldsymbol{x}^{(i)}) - y^{(i)})^2 \right] + \lambda \sum_{i=1}^{n} \boldsymbol{w}_i^2$$

L2 heavily penalizes larger parameters, preferring all smaller values.

L1-norm (Lasso Regression):

$$J(\boldsymbol{w}) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\boldsymbol{w}}(\boldsymbol{x}^{(i)}) - y^{(i)})^2 \right] + \lambda \sum_{i=1}^{n} |\boldsymbol{w}_i|$$

► L1 may set values of certain parameters to 0 → effectively "feature selection" (select the more important features, and zero out the rest).



Want to maximize the margin.



- 1. Find a point z on  $w \cdot x = 0$ . (Can always take z = 0)
- 2. Distance from  $\boldsymbol{x}$  to  $\boldsymbol{w} \cdot \boldsymbol{x} = 0$ =  $\frac{|(\boldsymbol{x} - \boldsymbol{z}) \cdot \boldsymbol{w}|}{||\boldsymbol{w}||}$ (Length of projection from  $\boldsymbol{x} - \boldsymbol{z}$  to  $\boldsymbol{w}$ )

Consider the decision boundary  $\boldsymbol{w} = \begin{vmatrix} 3 \\ 4 \end{vmatrix}$ . Given that the sample  $\mathbf{x} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$  is a support vector, what is the size of the margin? A. 2.2 B. 2.8 **Solution**. We have  $z = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ . C. 4.4 D. 5.6 Distance from x to decision boundary E. None of the above  $=\frac{1}{3^2+4^2}\left(\begin{bmatrix}2\\2\end{bmatrix}-\begin{bmatrix}0\\0\end{bmatrix}\right)\cdot\begin{bmatrix}3\\4\end{bmatrix}=\frac{6\cdot 3+8\cdot 4}{5}=2.8.$ The size of the margin is  $2.8 \times 2 = 5.6$ .

Primal Formulation:

minimize<sub>w</sub> 
$$rac{1}{2} \| oldsymbol{w} \|^2$$
  
subject to  $y^{(i)}(oldsymbol{w} \cdot oldsymbol{x}^{(i)}) \geq 1$ 

- **Convex** optimization problem with linear constraints.
- Slow for high-dimensional data (computing w · x<sup>(i)</sup> requires O(m) time, m = number of features).



Dual Formulation: First,

$$\begin{aligned} \text{maximize}_{\boldsymbol{\alpha}} & \sum_{i} \alpha^{(i)} - \frac{1}{2} \sum_{i} \sum_{j} \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \left( \boldsymbol{x}^{(i)} \cdot \boldsymbol{x}^{(j)} \right) \\ \text{subject to} & \sum_{i} \alpha^{(i)} y^{(i)} = 0 \text{ and } \alpha^{(i)} \ge 0 \end{aligned}$$

Then compute

$$oldsymbol{w} = \sum_i lpha^{(i)} y^{(i)} oldsymbol{x}^{(i)}$$

**Convex** optimization problem with linear constraints.

- **\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}** can be precomputed efficient when there are many features.
- For the **optimal solution**,  $\alpha^{(i)}$  for non-support vectors can always be 0.

(3,0)

X1

ector Machines  

$$y^{(1)}y^{(5)}(x^{(1)} \cdot x^{(5)}) = (-1)(1)\left(\begin{bmatrix} -2\\ -2 \end{bmatrix} \cdot \begin{bmatrix} 3\\ 0 \end{bmatrix}\right)$$
maximize<sub>\alpha</sub>  $\alpha^{(1)} + \alpha^{(2)} + \alpha^{(3)} + \alpha^{(4)} + \alpha^{(5)}$ 

$$-\frac{1}{2}(8\alpha^{(1)2} + 4\alpha^{(1)}\alpha^{(2)} + 4\alpha^{(1)}\alpha^{(3)} + 4\alpha^{(1)}\alpha^{(4)} + 6\alpha^{(1)}\alpha^{(5)}$$

$$+ 4\alpha^{(2)}\alpha^{(1)} + 4\alpha^{(2)2} + 4\alpha^{(2)}\alpha^{(3)} + 2\alpha^{(2)}\alpha^{(5)}$$

$$+ 4\alpha^{(3)}\alpha^{(1)} + 4\alpha^{(3)2} + 2\alpha^{(3)}\alpha^{(4)}$$

$$\begin{split} &+4\alpha^{(4)}\alpha^{(1)}+2\alpha^{(4)}\alpha^{(2)}+2\alpha^{(4)}\alpha^{(3)}+2\alpha^{(4)2}+3\alpha^{(4)}\alpha^{(5)}\\ &+6\alpha^{(5)}\alpha^{(1)}+6\alpha^{(5)}\alpha^{(2)}+3\alpha^{(5)}\alpha^{(4)}+9\alpha^{(5)2})\\ \text{subject to } &-\alpha^{(1)}-\alpha^{(2)}+\alpha^{(3)}+\alpha^{(4)}+\alpha^{(5)}=0 \text{ and } \alpha^{(i)}\geq 0 \end{split}$$

"Throws to an off-the-shelf solver" 
$$\alpha^{(1)} = 0, \alpha^{(2)} = 0.25, \alpha^{(3)} = 0.25, \alpha^{(4)} = 0, \alpha^{(5)} = 0$$

 $\overline{2}$ 

SVM Example

 $x_2$ 

(1, 1)

 $\times$ 

(0

(-2, 0)

× (−2,−2)

$$\therefore \boldsymbol{w} = -0 \begin{bmatrix} -2\\ -2 \end{bmatrix} - 0.25 \begin{bmatrix} -2\\ 0 \end{bmatrix} + 0.25 \begin{bmatrix} 0\\ 2 \end{bmatrix} + 0 \begin{bmatrix} 1\\ 1 \end{bmatrix} + 0 \begin{bmatrix} 3\\ 0 \end{bmatrix} = \begin{bmatrix} 0.5\\ 0.5 \end{bmatrix}$$

× (−2,−2)



 $\therefore$  The support vectors are (-2, 0), (0, 2) and (1, 1).



SVM Example

(b) Suppose we introduce another point,  $\mathbf{x}^{(6)}$ , with features  $\begin{bmatrix} -5 & 1 \end{bmatrix}^{\top}$  and label -1, then retrain the SVM. Will the learned model change?

Distance to decision boundary

$$=\frac{|\boldsymbol{x}^{(6)}\cdot\boldsymbol{w}|}{\|\boldsymbol{w}\|}=\frac{\left|\begin{bmatrix}-5\\1\end{bmatrix}\cdot\begin{bmatrix}0.5\\0.5\end{bmatrix}\right|}{\sqrt{0.5^2+0.5^2}}=\sqrt{2}\cdot 2.$$

Since  $x^{(6)}$  is not a support vector, the learned model will not change.



SVM Example

- (c) Let's remove data point  $x^{(2)}$  and retrain. What will happen to the model?
  - x<sup>(2)</sup> is a support vector, so the decision boundary is probably affected.
  - $\mathbf{x}^{(1)} = (-2, -2)$  becomes the new support vector.
  - ► The decision boundary should be shifted towards (-2, -2).



SVM Example

- (d) How would the results differ when we remove  $\mathbf{x}^{(3)}$  instead of  $\mathbf{x}^{(2)}$  and retrain the model?
- x<sup>(3)</sup> is a support vector, so the decision boundary is probably affected.
- The support vectors would be (-2, 0) and (1, 1).
- The decision boundary would be the perpendicular bisector of the line segment connecting them.

- (a) Write down the expression for the smallest distance of all points to a hyperplane defined by some  $\boldsymbol{w}$ .
  - ▶ Take *z* = **0**.
  - Distance of all points to a hyperplane =

$$=\frac{|(\boldsymbol{x}-\boldsymbol{z})\cdot\boldsymbol{w}|}{\|\boldsymbol{w}\||}=\frac{|\boldsymbol{x}\cdot\boldsymbol{w}|}{\|\boldsymbol{w}\|}.$$

- $\therefore$  The expression is  $\min_{i} \frac{|\boldsymbol{x} \cdot \boldsymbol{w}|}{\|\boldsymbol{w}\|}$ .
- (b) Write down the expression for maximising the smallest distance of all points to a hyperplane defined by some w.

• The expression is 
$$\max_{\boldsymbol{w}} \min_{i} \frac{|\boldsymbol{x} \cdot \boldsymbol{w}|}{\|\boldsymbol{w}\|}$$
.

- (c) Does this expression satisfy the correct classification constraints of the SVM, i.e., do the points lie on the correct side of the hyperplane?
  - No!

Our current optimization problem:

$$egin{array}{lll} \mathsf{maximize}_{oldsymbol{w}} & \min_i rac{|oldsymbol{x}\cdotoldsymbol{w}|}{\|oldsymbol{w}\|} \ \mathsf{subject to} & y^{(i)}(oldsymbol{w}\cdotoldsymbol{x}^{(i)}) > 0 \end{array}$$

- Observation 1:  $\boldsymbol{w} \cdot \boldsymbol{x}^{(i)} = \pm c$  for support vectors.
- Observation 2: We can scale w freely.
- ► ⇒ Fix  $\boldsymbol{w} \cdot \boldsymbol{x}^{(i)} = \pm 1$  by scaling  $\boldsymbol{w}$ .

New optimization problem:

$$\begin{aligned} \max_{i} \min_{i} \frac{|\boldsymbol{x} \cdot \boldsymbol{w}|}{\|\boldsymbol{w}\|} &= \frac{1}{\|\boldsymbol{w}\|} \\ \text{subject to } y^{(i)}(\boldsymbol{w} \cdot \boldsymbol{x}^{(i)}) \geq 1 \end{aligned}$$

• Maximize 
$$\frac{1}{\|\boldsymbol{w}\|} \Rightarrow$$
 Minimize  $\|\boldsymbol{w}\| \Rightarrow$  Minimize  $\frac{1}{2}\|\boldsymbol{w}\|^2$ 

Primal formulation:

$$\begin{array}{l} \text{minimize}_{\boldsymbol{w}} \ \frac{1}{2} \|\boldsymbol{w}\|^2\\ \text{subject to } y^{(i)}(\boldsymbol{w} \cdot \boldsymbol{x}^{(i)}) \geq 1 \end{array}$$



Hypothesis class:  $h_w(x) = w_0 x + b$ 

**Bias** measures how much the **expected** *predicted value* differs from the *actual value*. Intuition: The error of the **best model** when you are given **infinite amount of training data**.





Variance measures the consistency of predictions due to sampling.



Poll:

 If we increase the amount of training samples: The bias <u>remains unchanged</u> and the variance <u>decreases</u>.
 If we increase the degree of the polynomial in our hypothesis class: The bias <u>decreases</u> and the variance <u>increases</u>.



#### More samples (Lower variance)





#### Higher degree (Low bias; High variance)



#### **Bias-Variance Tradeoff**



#### **Bias-Variance Tradeoff**

Underfitting



Just Right



Overfitting



#### Q4. Bias and Variance

The model hypotheses are as below:

1. 
$$H_w(x) = w_0 + w_1 x$$
  
2.  $H_w(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_{10} x^{10}$ 



Main Idea:

- Both bias and variance contribute to the mean squared error.
- Bias does not decrease with the number of training samples. Bias causes error in both the training set and the test set.
- Variance decreases with the number of training samples. Variance causes a difference between the error in the training set and that in the test set.

#### Q4. Bias and Variance

The model hypotheses are as below:

1.  $H_w(x) = w_0 + w_1 x$  High bias

2.  $H_w(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_{10} x^{10}$  High variance



# Recap: Kernel Trick

Using the dual formulation:

$$\begin{aligned} & \text{maximize}_{\alpha} \ \sum_{i} \alpha^{(i)} - \frac{1}{2} \sum_{i} \sum_{j} \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \left( \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} \right) \\ & \text{subject to} \ \sum_{i} \alpha^{(i)} y^{(i)} = 0 \text{ and } \alpha^{(i)} \ge 0 \end{aligned}$$

We transform the features in  $\mathbf{x}^{(i)}$  to  $\phi(\mathbf{x}^{(i)})$ . How does that impact the training?

- ▶ Not much, we have to change  $\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$  to  $\phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}^{(j)})$ .
- We define a kernel function K(x<sup>(i)</sup>, x<sup>(j)</sup>) = φ(x<sup>(i)</sup>) ⋅ φ(x<sup>(j)</sup>) to compute this efficiently.

**Important Property of a Kernel Function**  $K(u, v) = \phi(u) \cdot \phi(v)$  for some function  $\phi$ 

#### Recap: Kernel Trick Example: $K(u, v) = (u \cdot v)^2$ , where $u, v \in \mathbb{R}^3$ .

K(

$$\begin{aligned} \mathbf{u}, \mathbf{v}) &= (\mathbf{u} \cdot \mathbf{v})^2 \\ &= \left(\sum_i u_i v_i\right)^2 \\ &= u_1^2 v_1^2 + u_2^2 v_2^2 + u_3^2 v_3^2 + 2u_1 u_2 v_1 v_2 + 2u_1 u_3 v_1 v_3 + 2u_2 u_3 v_2 v_3 \\ &= \begin{bmatrix} u_1^2 \\ u_2^2 \\ u_3^2 \\ \sqrt{2}u_1 u_2 \\ \sqrt{2}u_1 u_3 \\ \sqrt{2}u_2 u_3 \end{bmatrix} \cdot \begin{bmatrix} v_1^2 \\ v_2^2 \\ v_3^2 \\ \sqrt{2}v_1 v_2 \\ \sqrt{2}v_1 v_3 \\ \sqrt{2}v_2 v_3 \end{bmatrix} \\ &= \phi(\mathbf{u}) \cdot \phi(\mathbf{v}) \end{aligned}$$

# Recap: Kernel Trick

Dual formulation with kernel trick:

maximize 
$$\sum_{i} \alpha^{(i)} - \frac{1}{2} \sum_{i} \sum_{j} \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \times \mathcal{K}(\boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)})$$
subject to 
$$\sum_{i} \alpha^{(i)} y^{(i)} = 0 \text{ and } \alpha^{(i)} \ge 0$$

Does not slow down training, but "implicitly" considers all transformed features!

#### Bonus. Gaussian Kernel

Prove that the Gaussian Kernel,  $K(\boldsymbol{u}, \boldsymbol{v}) = e^{-\frac{\|\boldsymbol{u}-\boldsymbol{v}\|^2}{2\sigma^2}}$ , has infinite dimensional features. (We assume  $\sigma^2 = 1$  for simplicity.)

Extra Slide