### CS3244 : Machine Learning

Semester 2 2023/24

### Lecture 4 : Linear Models and Support Vector Machine

#### Xavier Bresson

https://twitter.com/xbresson



Department of Computer Science National University of Singapore (NUS)



# Outline

#### • Admin

- Three applications of linear models
- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

# Outline

#### • Admin

#### • Three applications of linear models

- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

## Three learning problems

• What are the simplest models to solve these learning tasks?



## Three learning problems

- What are the simplest models to solve these learning tasks?
  - Linear models baseline to any machine learning project.



### Linear models

- Hypothesis space H, i.e. the space of all solutions, of linear models :
  - Straight lines (2D space), planes (3D space) and hyper-planes (>3D spaces).
- Parameters of linear models are the slopes and the bias.



55

# Outline

#### • Admin

• Three applications of linear models

#### • Classification

- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

## Three learning problems



## Linear model for credit approval

• Goal : Design a linear model to approve or not a credit to an individual based on personal information or features.



## Example

- Extend credit using applicant's age
  - If age is higher than 25 years, accept credit.
  - Otherwise, reject.
- Linear model :  $f_{\theta}(\mathbf{x}) = x_1 25 > 0 \Rightarrow Approve$  $< 0 \Rightarrow Reject$ With  $f_{\theta}(\mathbf{x}) = \mathbf{\theta}^{\top} \mathbf{x} = \theta_1 x_1 + \theta_0$ ,  $\mathbf{\theta} = (\theta_1, \theta_0) = (1, -25)$  and  $\mathbf{x} = (x_1, 1)$

• Example :

$$f_{\theta}(\mathbf{x}) = 32 - 25 = 7 > 0 \Rightarrow Approve$$

### Binary classification

• Binarize a linear function :

For a data input  $\mathbf{x} = (x_1, x_2, ..., x_n)$  (attributes of a customer)

Approve credit if:  $\mathbf{\theta}^{\top} \mathbf{x} = \sum_{i=1}^{n} \theta_i x_i > 0$ Deny credit otherwise:  $\mathbf{\theta}^{\top} \mathbf{x} = \sum_{i=1}^{n} \theta_i x_i < 0$ 

The linear formula f can be written as:



### Decision boundary in 2D



# Decision boundary in 3D



# Outline

#### • Admin

- Three applications of linear models
- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

# Three learning problems



### Linear model for credit amount

• Goal : Design a linear model to estimate the credit amount given to an individual based on personal information or features.

Criterion	Value	
Age	32 years	<i>x</i> <sub>1</sub>
Salary	40 K	$x_2$
Debt	26 K	
Years in Job	1 year	
Years at Current Residence	3 years	x <sub>d</sub>



### Loss function

- Quality of a predictive function for a training set is evaluated with a loss function.
  - Training set : Collected/historical data  $(\mathbf{X}, \mathbf{y}) = (\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$ where **X** is a n x d data matrix and **y** is a m-dim label vector
    - For classification,  $y^{(j)} = \pm 1$  and for regression  $y^{(j)} \in \mathbb{R}$ .
  - Loss function computes the fit of the predictive function  $f_{\theta}(x)$  to the label y:



How good is the prediction

## Illustration

- The training set is composed of n data points, each data feature is 1-dimensional, i.e.  $x = x_1$  and the label is y.
- Find the linear model  $\theta$  that fits all the data points as best as possible.



Hypothesis function :  $f_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1$ 

Loss value:  $L(\theta_{green}) \ge L(\theta_{blue})$ 

### Illustration

• The best possible linear model  $f_{\theta}(\mathbf{x})$  will have the minimum loss value of  $L(\theta)$ .



Hypothesis function :  $f_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1$ 

 $f_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ 

 $\min_{\theta} L(\theta) = \frac{1}{n} \sum_{j=1}^{n} (f_{\theta}(\mathbf{x}^{(j)}) - y^{(j)})^{2}$ 

# Outline

#### • Admin

- Three applications of linear models
- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

### Parameter estimation

• How to compute the best parameters  $\theta$  that minimizes the MSE loss function?

$$\min_{\theta} L(\theta) = \frac{1}{n} \sum_{j=1}^{n} (f_{\theta}(\mathbf{x}^{(j)}) - y^{(j)})^2$$

- Two approaches
  - Normal equations
  - Gradient descent

### Normal equations

$$L = \frac{1}{n} \sum_{j=1}^{n} (\boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}^{(j)} - \mathbf{y}^{(j)})^2$$
$$= \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2$$

where 
$$\mathbf{X} = \begin{bmatrix} -\mathbf{x}^{(1)^{\mathsf{T}}} - \\ -\mathbf{x}^{(2)^{\mathsf{T}}} - \\ -\mathbf{x}^{(3)^{\mathsf{T}}} - \\ \vdots \\ -\mathbf{x}^{(n)^{\mathsf{T}}} - \end{bmatrix}$$
,  $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ \vdots \\ y^{(n)} \end{bmatrix}$  and  $\mathbf{\theta} = \begin{bmatrix} \theta^{(1)} \\ \theta^{(2)} \\ \theta^{(3)} \\ \vdots \\ \theta^{(d)} \end{bmatrix}$ 

# Normal equations

• Example of  $X\theta$  and y:

$$\mathbf{X}\mathbf{\Theta} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} \end{bmatrix} \times \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} = \begin{bmatrix} \theta_1 x_1^{(1)} + \theta_2 x_2^{(1)} + \theta_3 x_3^{(1)} \\ \theta_1 x_1^{(1)} + \theta_2 x_2^{(2)} + \theta_3 x_3^{(3)} \end{bmatrix} = \begin{bmatrix} \mathbf{\Theta}^\top \mathbf{x}^{(1)} \\ \mathbf{\Theta}^\top \mathbf{x}^{(2)} \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \end{bmatrix}$$

### Normal equations

• Expand L :

$$L = \frac{1}{n} \| \mathbf{X} \boldsymbol{\theta} - \mathbf{y} \|^{2}$$
$$= \frac{1}{n} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y}) = \frac{1}{n} ((\mathbf{X} \boldsymbol{\theta})^{\top} - \mathbf{y}^{\top}) (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$$
$$= \frac{1}{n} ((\mathbf{X} \boldsymbol{\theta})^{\top} \mathbf{X} \boldsymbol{\theta} - (\mathbf{X} \boldsymbol{\theta})^{\top} \mathbf{y} - \mathbf{y}^{\top} \mathbf{X} \boldsymbol{\theta} + \mathbf{y}^{\top} \mathbf{y})$$
$$= \frac{1}{n} (\boldsymbol{\theta}^{\top} \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} - 2\boldsymbol{\theta}^{\top} \mathbf{X}^{\top} \mathbf{y} + \mathbf{y}^{\top} \mathbf{y}), \text{ with } \boldsymbol{\theta}^{\top} \mathbf{X}^{\top} \mathbf{y} = (\mathbf{X} \boldsymbol{\theta})^{\top} \mathbf{y} = \mathbf{y}^{\top} \mathbf{X} \boldsymbol{\theta}$$

• Compute gradient (with calculus) and get solution :

$$\frac{\partial \mathbf{L}}{\partial \boldsymbol{\theta}} = \frac{1}{n} (2 \cdot \mathbf{X}^{\mathsf{T}} \mathbf{X} \boldsymbol{\theta} - 2 \cdot \mathbf{X}^{\mathsf{T}} \mathbf{y}) = \mathbf{0}$$
  

$$\Rightarrow \mathbf{X}^{\mathsf{T}} \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^{\mathsf{T}} \mathbf{y} \Rightarrow \boldsymbol{\theta} = ((\mathbf{X}^{\mathsf{T}} \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}}) \mathbf{y}$$
  
Pseudo-inverse of **X**  
Computationally expensive

re to compute if d is large.

# Outline

#### • Admin

- Three applications of linear models
- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

## Three learning problems



- Linear classification predicts "hard" classes, e.g.  $\pm 1$ .
- Linear regression predicts scalar target, e.g.  $\in \mathbb{R}$ .
- Logistic regression predicts "soft" classes, i.e.  $f_{\theta}(\mathbf{x})$  is interpreted as a probability of class  $\in [0,1]$ .
- Example
  - Prediction of a heart attack
  - Input  $\mathbf{x}$ : Cholesterol level, age, diabetic, etc.
  - Output  $f_{\theta}(\mathbf{x})$ : Likelihood of a heart attack  $\in [0,1]$
  - $\theta^{\top} \mathbf{x} \in \mathbf{R}$  is the risk score.

- Training data :  $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$  with
  - $\mathbf{x}^{(j)}$ : Person's health information, i.e. data features
  - $y^{(j)}$ : Person has a heart attack or not
- Labels y are binary, i.e.  $y^{(j)} = \pm 1$ .
  - We know with certainty whether a person got a heart attack,  $y^{(j)} = +1$ , or not,  $y^{(j)} = -1$ .
- Prediction function  $f_{\theta}(\mathbf{x})$  is the probability in [0,1] that the person had a heart attack.
  - Hence, the probability of a person with no heart attack is 1  $f_{\theta}(\mathbf{x})$ .

$$P_{\theta}(y|\mathbf{x}) = \begin{cases} f_{\theta}(\mathbf{x}) & \text{for } y = +1\\ 1 - f_{\theta}(\mathbf{x}) & \text{for } y = -1 \end{cases}$$

Probability of having class y given the input data x Conditioned probability a.k.a. likelihood probability

• Prediction  $f_{\theta}(x)$  can be designed as :

 $f_{\theta}(\mathbf{x}) = \text{sigmoid}(\boldsymbol{\theta}^{\top}\mathbf{x}) \in [0,1]$ 

- Choice of logistic function :
  - Function sigmoid(.)=g(.) is a smooth step function or soft thresholding function.

$$g(s) = \frac{\exp(s)}{1 + \exp(s)} = \frac{1}{1 + \exp(-s)}$$

$$g(-s) = 1 - g(s)$$

$$0 \qquad g(s)$$

$$g(-s) = 1 - g(s)$$

$$0 \qquad g(s)$$

$$g(s)$$

$$g(s)$$

$$g(s)$$

- Loss function for logistic regression
  - Estimate how good is the predictive function  $f_{\theta}(\mathbf{x}) = \text{sigmoid}(\boldsymbol{\theta}^{\top}\mathbf{x})$  to infer correctly the class y of x :

$$\min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \frac{1}{n} \sum_{j=1}^{n} \log(1 + \exp(-y^{(j)} \boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})) \qquad \begin{array}{c} \text{Logistic r} \\ \text{a.k.a cr} \end{array}$$
Number of training data

Logistic regression function a.k.a cross-entropy loss

Explained in the next slides

• Understanding the logistic regression loss for one data point (x,y):

Sigmoid  

$$P_{\theta}(y|\mathbf{x}) = \begin{cases} g(\theta^{\mathsf{T}}\mathbf{x}) = f_{\theta}(\mathbf{x}) & \text{for } y = +1 \\ g(-\theta^{\mathsf{T}}\mathbf{x}) = 1 - g(\theta^{\mathsf{T}}\mathbf{x}) = 1 - f_{\theta}(\mathbf{x}) & \text{for } y = -1 \\ g(-s) = 1 - g(s) \end{cases} \quad \text{for } y = -1 = g(y.\theta^{\mathsf{T}}\mathbf{x})$$

$$P_{\theta}(y|\mathbf{x}) = g(y.\theta^{\mathsf{T}}\mathbf{x}) = \frac{1}{1 + \exp(-y.\theta^{\mathsf{T}}\mathbf{x})} \quad \text{with } g(s) = \frac{1}{1 + \exp(-s)}$$

$$\log P_{\theta}(y|\mathbf{x}) = -\log(1 + \exp(-y.\theta^{\mathsf{T}}\mathbf{x}))$$

$$\max_{\theta} P_{\theta}(y|\mathbf{x}) \Leftrightarrow \max_{\theta} \log P_{\theta}(y|\mathbf{x}) = \max_{\theta} - \log(1 + \exp(-y.\theta^{\mathsf{T}}\mathbf{x}))$$

$$\Leftrightarrow \min_{\theta} \log(1 + \exp(-y.\theta^{\mathsf{T}}\mathbf{x}))$$
Monotonous
increasing function

- Understanding the logistic regression loss for all data points  $\mathbf{X} = (\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})$ :
  - Suppose that the data points are i.i.d. samples (independent and identically distributed).

$$P(y^{(1)}, ..., y^{(n)} | \mathbf{x}^{(1)}, ..., \mathbf{x}^{(n)}) = \prod_{j=1}^{n} P(y^{(j)} | \mathbf{x}^{(j)})$$
  
$$\log P(y^{(1)}, ..., y^{(n)} | \mathbf{x}^{(1)}, ..., \mathbf{x}^{(n)}) = \log \prod_{j=1}^{n} P(y^{(j)} | \mathbf{x}^{(j)}) = \sum_{j=1}^{n} \log P(y^{(j)} | \mathbf{x}^{(j)})$$
  
$$\min_{\theta} L(\theta) = \frac{1}{n} \sum_{j=1}^{n} \log(1 + \exp(-y^{(j)} \theta^{\top} \mathbf{x}^{(j)})) \qquad \log \prod_{j=1}^{n} = \sum_{j=1}^{n} \log(1 + \exp(-y^{(j)} \theta^{\top} \mathbf{x}^{(j)}))$$

# Outline

#### • Admin

- Three applications of linear models
- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

## Gradient optimization

- Gradient optimization algorithm (simplest optimization algorithm)
  - Work well in high-dimensional spaces (because convergence is independent of data dimensionality)
- Two versions : Gradient descent/ascent
  - Start randomly
  - Move in the direction of the steepest descent/ascent.
  - Stop when reaching the minimum/maximum.



- Predict amount of credit y with data feature x (single feature) : ٩
  - Prediction function :  $f_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 \mathbf{x}$ 0
  - $\theta_0, \theta_1$   $L(\theta = (\theta_0, \theta_1)) = \frac{1}{n} \sum_{i=1}^n (f_\theta(\mathbf{x}^{(i)}) \mathbf{y}^{(i)})^2$ Parameters : 0 Loss function : 0
  - Loss minimized by gradient descent (introduced next slides): 0



#### $\min_{\theta} L(\theta = (\theta_0, \theta_1))$

- Univariate linear regression
  - (Single) parameter  $\theta$  : slope of the line
- Suppose loss function is a convex function of  $\theta$ , e.g.  $L(\theta) = (\theta x y)^2$



- Multivariate linear regression
  - Two parameters  $\theta_0$ ,  $\theta_1$ : bias and slope of the line
- Suppose loss function is a convex function of  $\theta = (\theta_0, \theta_1)$ , e.g.  $L(\theta_0, \theta_1) = (\theta_0 + \theta_1 x y)^2$



- $\theta(t)$ : Vector of all parameters at iteration t
- Start at  $\boldsymbol{\theta}(t=0) = random$
- Iterative method
  - Update the value of parameters :  $\boldsymbol{\theta}(t+1) = \boldsymbol{\theta}(t) + \boldsymbol{v}$
  - What is the best direction  $\mathbf{v}$  ?
    - Direction **v** must minimize as much as possible the loss function L.
    - **v** = Opposite direction of the steepest slope (justification given next slides)



• Direction  $\mathbf{v}$  of the steepest descent :



 $\mathbf{v} = -\eta \frac{\partial L}{\partial \theta} \int_{0}^{\infty} \frac{\partial L}{\partial \theta}$  $\theta(t)$  $\theta(t)$  $\mathbf{v}(t) = \Delta \theta(t) = -\eta \frac{\partial L}{\partial \theta}(\theta(t))$ 

 $L(\theta)$ 

- Step size  $\eta$  controls how fast the gradient descent algorithm descends to the minimum.
- Finding the right value of the step size is heuristic :



- Termination condition / when to stop the iterative scheme ?
  - Natural choice : when gradient < (arbitrary) threshold
- However, lots of saddle points / flat regions in high-dimensional spaces.





# Stochastic gradient descent

- A fast variation of gradient descent that considers only a small set of data points to update the value of the parameters is stochastic gradient descent (SGD):
  - Pick a (random) small subset of m training data (e.g. a single/512 data points), i.e.  $(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$
  - Compute the loss value for this subset, i.e.  $L = \frac{1}{m} \sum_{k=1}^{m} L(f_{\theta}(\mathbf{x}^{(k)}), \mathbf{y}^{(k)})$
  - Compute the gradient for this subset, i.e.  $\mathbf{v} = \frac{1}{m} \sum_{k=1}^{m} -\nabla L(f_{\theta}(\mathbf{x}^{(k)}), \mathbf{y}^{(k)})$
  - Update parameters, i.e.  $\theta(t+1) = \theta(t) \eta v$
- Advantages
  - Faster update of parameters : Gradient computed with m=512 rather than all data points.
  - Stochastic optimization : Helps escape saddle points in high-dimensional space
  - Simple to implement



## Stochastic gradient descent

• GD vs. SGD with by a mini-batch of m=10 samples



## Summary

- Two methods for solving the linear regression task :
  - Gradient Descent
    - Works well for very large n, #training data, with SGD (300B tokens with GPT3)
    - Works well for very large d, #data features (d=1M with 1,000×1,000 images)
    - Works well for very large  $|\theta|$ , #parameters features (175B with GPT3)
    - Very slow and requires to select time step  $\eta$
  - Normal Equations
    - Very fast for d = O(10<sup>3</sup>) data points, do not require to choose  $\eta$
    - Need to compute  $(\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}$ ,  $O(d^3)$  operation but faster approximations exist
- Gradient Descent is a universal optimization technique as long as the considered loss is continuous and differentiable (as gradient is required).
  - GD does not work for discrete losses like win/lose at the game of Go.

# Outline

#### • Admin

- Three applications of linear models
- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

- Consider the binary classification task and linear functions as hypothesis space.
- The goal is to find a linear separator that partitions the feature space into two regions.
- Generally, there exist several possible linear separators.
- How to select an "optimal" linear separator, and how to define "optimal"?



- SVM technique (1964) aims at maximizing the margin between the two classes.
- Maximizing the margin provides better results.



- Hyper-plane equation : {  $\mathbf{x}$  such that  $f_{\theta}(\mathbf{x}) = \mathbf{\theta}^{\top} \mathbf{x} + \mathbf{b} = \mathbf{cte}$  }
- Hyper-planes  $f_{\theta}(\mathbf{x}) = \pm 1$  are margin planes.
- Hyper-plane  $f_{\theta}(\mathbf{x}) = 0$  is the class separator.
- Training data points +1 lie above the hyperplane and -1 data points below.
- Parameters  $\boldsymbol{\theta}$  controls the orientation/slope of the plane, and b is the offset/bias.
- The distance *d* between the two margin planes is computed to be as large as possible.
- What are the parameters  $\theta$  that maximizes the distance d between the two margins? (response next slide)



- What are the parameters  $\boldsymbol{\theta}$  that maximizes the margin?
- Optimal value  $\theta$  is solution of a constrained quadratic optimization problem :

Given  $f_{\theta}(\mathbf{x}) = \mathbf{\theta}^{\top} \mathbf{x} + \mathbf{b}$ with  $f_{\theta}(\mathbf{x}_{+1}) = \mathbf{\theta}^{\top} \mathbf{x}_{+1} + \mathbf{b} = +1$  $f_{\theta}(\mathbf{x}_{-1}) = \mathbf{\theta}^{\top} \mathbf{x}_{-1} + \mathbf{b} = -1$ We have  $f_{\theta}(\mathbf{x}_{+1}) - f_{\theta}(\mathbf{x}_{-1}) = \theta^{T}(\mathbf{x}_{+1} - \mathbf{x}_{-1}) = +2$ If we define the margin vector as  $\mathbf{d} = \mathbf{x}_{+1} - \mathbf{x}_{-1}$  then  $\boldsymbol{\theta}^{\mathsf{T}}(\mathbf{x}_{+1} - \mathbf{x}_{-1}) = +2 \Rightarrow \|\boldsymbol{\theta}\|. \|\boldsymbol{d}\| = +2 \Rightarrow \mathbf{d} = \frac{2}{\|\boldsymbol{\theta}\|}$  $\max_{\mathsf{d}} \mathbf{d} = \frac{2}{\|\boldsymbol{\theta}\|} \Leftrightarrow \min_{\boldsymbol{\theta}} \|\boldsymbol{\theta}\|^{2} \text{ such that}$  $\begin{cases} \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}) \geq +1 \text{ for } \mathbf{x} \in \mathbf{R}_{+1} \\ \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}) \leq -1 \text{ for } \mathbf{x} \in \mathbf{R}_{-1} \end{cases}$ Length of the margin vector



Only a subset of the data points have value  $f(\mathbf{x}) = \pm 1$ . These are called support vectors (SV).

# Outline

#### • Admin

- Three applications of linear models
- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

- Real-world data is noisy.
- In case of data non-linearly separable, i.e. with outliers, there is no mathematical solution to standard/hard-margin SVM.
  - In other words, there exists no linear separator that can split the two classes perfectly, i.e. without errors.



- SVM can be improved to deal with outliers.
  - Soft-margin SVM technique, 1995
  - Idea is to introduce a slack variable  $e^{(j)}$  for each data point that represents the prediction error.
  - These errors  $e^{(j)}$  will be minimized while simultaneously maximizing the margin :

• Effect of varying C, the regularization parameter.

Goal is to find the largest margin. When C is small, more misclassification errors are allowed.

Note that the margin is large then.



When C is large, less misclassification error are allowed, possibly none. Note that the margin is small then.

- Soft-margin SVM penalizes
  - Misclassifications,
  - Correct classifications that fall inside the margin.
- Loss of soft-margin : Hinge loss (popular loss function)

$$\begin{split} \min_{\boldsymbol{\theta}, \mathbf{e}} \|\boldsymbol{\theta}\|^{2} + C. \sum_{j=1}^{n} e^{(j)} \text{ such that } \begin{cases} f_{\boldsymbol{\theta}}(\mathbf{x}^{(j)}) \geq +1 - e^{(j)} \text{ for } \mathbf{x}^{(j)} \in R_{+1} \\ f_{\boldsymbol{\theta}}(\mathbf{x}^{(j)}) \leq -1 + e^{(j)} \text{ for } \mathbf{x}^{(j)} \in R_{-1} \\ e^{(j)} \geq 0, \ C \geq 0 \end{cases} \overset{d \leq 1} \underbrace{d \geq 1}_{L_{\text{hin}} > 0} \underbrace{d \geq 1}_{1 \ L_{\text{Hin}} = 0} \\ \dim_{\boldsymbol{\theta}} \|\boldsymbol{\theta}\|^{2} + C. \sum_{j=1}^{n} \max(0, 1 - y^{(j)}. f_{\boldsymbol{\theta}}(\mathbf{x}^{(j)})) \underbrace{L_{\text{Hin}} (d^{(j)} = y^{(j)}. f_{\boldsymbol{\theta}}(\mathbf{x}^{(j)}))}_{(\text{Hinge loss})} \overset{d \geq 1}{(\text{correctly classified})} \overset{d \geq 1}{f_{\boldsymbol{\theta}}(\mathbf{x}^{(j)}) \leq -1 \text{ and } y^{(j)} = -1} \end{aligned}$$

[Optional] Proof given in <u>http://image.diku.dk/imagecanon/material/cortes\_vapnik95.pdf</u>

Xavier Bresson

 $\mathrm{L}_{\mathrm{Hin}}$ 

.

# Outline

#### • Admin

- Three applications of linear models
- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

# Kernel SVM

- Linear models are limited to linearly separable data points.
- How to classify complex/non-linear datasets with linear separators?
- Idea is to map the data from their original space  $\mathbb{R}^d$ , where classes can only be separated with non-linear functions, to a new higherdimensional space  $\mathbb{R}^b$ ,  $b \gg d$ , where classes can be distinguished with linear functions :  $\mathbf{x} = (x_1, x_2, ..., x_n) \rightarrow \mathbf{z} = (z_1, z_2, ..., z_n)$





# Kernel SVM



# Kernel SVM

#### • Kernel trick

- Processing data points  $\Phi(\mathbf{x})$  in the higher-dimensional space  $R^b$  is time and memory consuming. Kernel trick avoids computing  $\Phi(\mathbf{x}) \in R^b$ .
- Instead, we will compute kernel value  $K(\mathbf{x},\mathbf{y})$  in  $\mathbb{R}^d$  with  $d\ll b$ .
- Different kernels exist such as polynomial/Gaussian kernels :  $K(x,y)=(1+x^Ty)^p$ ,  $K(x,y)=\exp(-x^Ty/\sigma)$
- Kernel SVM solutions are computed by solving a quadratic optimization problem :

$$f_{\theta}(\mathbf{x}) = \mathbf{\theta}^{\top} \mathbf{x} \stackrel{\Phi}{\rightarrow} f_{\theta}(\mathbf{x}) = \mathbf{\theta}^{\top} \Phi(\mathbf{x})$$
$$\mathbf{\theta} = \sum_{j=1}^{n} \alpha^{(j)} y^{(j)} \mathbf{x}^{(j)} \stackrel{\Phi}{\rightarrow} \mathbf{\theta} = \sum_{j=1}^{n} \alpha^{(j)} y^{(j)} \Phi(\mathbf{x}^{(j)})$$
Primal variable Dual variable

$$\mathbf{f}_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{\theta}^{\top} \boldsymbol{\Phi}(\boldsymbol{x}) = \sum_{j=1}^{n} \boldsymbol{\alpha}^{(j)} \boldsymbol{y}^{(j)} \boldsymbol{\Phi}(\boldsymbol{x}^{(j)})^{\mathsf{T}} \boldsymbol{\Phi}(\boldsymbol{x})$$

with  $\min_{0 \le \alpha \le C} \alpha^\top Q \alpha - \alpha^\top 1$  such that  $\alpha^\top y = 0$ , Q = YKY, Y = diag(y),  $K(x,y) = \Phi(x)^T \Phi(y)$ [proof not included -- it will not be assessed] [Optional] Proof given in <u>http://image.diku.dk/imagecanon/material/cortes\_vapnik95.pdf</u>

Xavier Bresson

58

# Outline

#### • Admin

- Three applications of linear models
- Classification
- Regression
- Normal equations
- Logistic regression
- Gradient descent
- Support Vector Machine
- Soft-margin SVM
- Kernel SVM
- Conclusion

## Convention

- In machine learning, linear models are linear functions w.r.t. data features x, not the parameters  $\theta$ .
- Remark 1: Mathematically, a linear function is a function that satisfies the additivity property, i.e. f(x+y) = f(x)+f(y). For example,  $f_{\theta}(x)=\theta_1x_1+\theta_2x_2$  is a linear function w.r.t.  $x = (x_1, x_2)$ .
- However, the affine function  $f_{\theta}(x) = \theta_1 x_1 + \theta_2 x_2 + \theta_0$  is not a linear function, but we make an abuse of notation and we still refer to the affine function as a "linear" model, but it is \*not\* true mathematically.
- Remark 2: Another abuse of notation is to disguise the affine function as a linear function as follows:  $f_{\theta}(x) = \theta_1 x_1 + \theta_0 = \theta^T x$  with  $\theta = (\theta_1, \theta_0)$  and  $x = (x_1, 1)$ . So it seems like a linear function w.r.t. x but it is not because it does not satisfy the additivity property, i.e. f(x+y) = f(x) + f(y).
- The same abuse can be used to pretend to represent quadratic function as "linear" function, i.e.  $f_{\theta}(x) = \theta_1 x + \theta_2 x^2 + \theta_0 = \theta^T x$ , where  $\theta = (\theta_2, \theta_1, \theta_0)$  and  $x = (x^2, x, 1)$ . Again, the representation seems linear w.r.t. x but mathematically it is not.

## Conclusion

- Linear models can be used for classification and regression tasks.
- Very well-established techniques, which fit optimally CPU/GPU acceleration hardware.
  - BLAS/LAPACK for CPU and CUDA for GPU
- But very limited expressivity, only perform well for linearly separable data points.
- Kernel SVM enhances their expressivity with non-linear separators but requires to handcraft a kernel operator.
  - Deep learning has significantly surpassed SVM techniques.

### References

- Prof Min-Yen Kan, CS3244 NUS, Machine Learning, 2022
  - <u>https://knmnyn.github.io/cs3244-2210</u>
- Prof Xavier Bresson, CS6208 NUS, Advanced Topics in Artificial Intelligence, 2023

