

### §3 compute the gradients

→ The gradients of the cost function with respect to  $m$  and  $c$  are computed as follows:

① gradient with respect to  $m$

$$\frac{\partial J}{\partial m} = \frac{1}{N} \sum_{i=1}^N -x_i (y_i - (mx_i + c))$$

②  $g$  with respect to  $c$

$$\frac{\partial J}{\partial c} = \frac{1}{N} \sum_{i=1}^N x_i (y_i - (mx_i + c))$$

### §3 update the parameters

→ using the learning rate  $\alpha$ , update the parameters

- modify  $m$  &  $c$

$$m = m - \alpha \frac{\partial J}{\partial m}$$

$$c = c - \alpha \frac{\partial J}{\partial c}$$

§4 iterate through the above steps for a certain number of iterations or until the cost function converges to a minimum.

§5 The process stops when the cost function  $J(m, c)$  doesn't change significantly b/w iterations.

- ② → linear regression model trained on 500 houses  
→ model predicts house prices based on square  
→ model performs poorly on new data
- ⓐ possible reasons for poor performance
- ① overfitting → the model performs well on training data but poorly on new data
- ② underfitting → the model is too simple and does not capture the patterns in the data
- ③ Data quality issues : noisy or inconsistent data, missing values, or incorrect labels
- ④ feature scaling → poorly scaled features
- ⑤ outliers → extreme values in the data can skew the model
- ⑥ insufficient data → 500 data points may not be enough to train a robust model

## ⑥ Improve their model's performance

- use cross-validation to check model generalization
- feature engineering: add relevant features
- use data augmentation or increase dataset size
- feature scaling: Normalize or standardize the square footage variable to improve gradient descent efficiency
- Handle data issues: clean data by removing outliers, filling missing values and ensuring consistency.
- Try diff models: Df, rf or nn for better performance.

③

polynomial regression eqn for two features with degree 2.

$$y = w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2$$

parameters to be estimated

$w_0 \rightarrow$  intercept of function

$w_1 \rightarrow$  coefficient for  $x_1$

$w_2 \rightarrow$  coefficient for  $x_2$

$w_3 \rightarrow$  coefficient for  $x_3$

$w_4 \rightarrow$  coefficient for  $x_4$

$w_5 \rightarrow$  coefficient for  $x_5$

→ These parameters are estimated using gradient descent or the Normal equation

$$\theta = (x^T x)^{-1} x^T y$$

④

→ Overfitting

↳ the model learns noise & pattern specific to the training data.

Symptoms:

① very high accuracy on training but poor performance on test data

② the model might memorize data points rather than generalizing.

### Example

- ① suppose the dataset contains 500 houses, and the model includes highly specific features like street name, exact house age, or very high-degree polynomial terms.

- ② The model might learn patterns that only apply to the training set, but this rule won't generalize to new houses.

Result 1. High training accuracy  $\rightarrow 98\%$

low test accuracy  $\rightarrow$  less 60%

### Under fitting

- $\rightarrow$  The model is too simple to capture important patterns in the data.

### Sym

- $\rightarrow$  poor performance on both training & test data  
 $\rightarrow$  High bias (unable to capture relationship properly)

-  ① suppose the model only considers square foot -ge but ignores location & no. of rooms.

- ② It assumes house prices are strictly linear (e.g. price =  $200 \times \text{sqf}$ ), even though prices also depend on factors like house age, proximity to school etc.

Result, low accuracy on both training & test

## → How to improve the model

### → Overfitting

#### ① Reduce Complexity

- ↳ reduce the degree of polynomial regression
- ↳ remove irrelevant features like exact street names.

#### ② Increase Training data

#### ③ split data into multiple training/test subsets to ensure the model generalizes well.

### → Underfitting

#### ① add more features

- ↳ include relevant variables like location, crime rate

#### ② Use a more flexible model

- ↳ try DT, RF, or NN

#### ③ use feature engineering

## Cross-validation

### ③ Hyperparameter tuning

↳ Hyperparameter tuning is the process of finding the best combination of hyperparameters to optimize a model's performance.

→ In polynomial regression, key hyperparameters

includes

① degree of the polynomial

② Regularization strength (prevents overfitting)

param\_grid = {

'poly-degree': [1, 2, 3],

'ridge-alpha': [0.01, 0.1, 1]

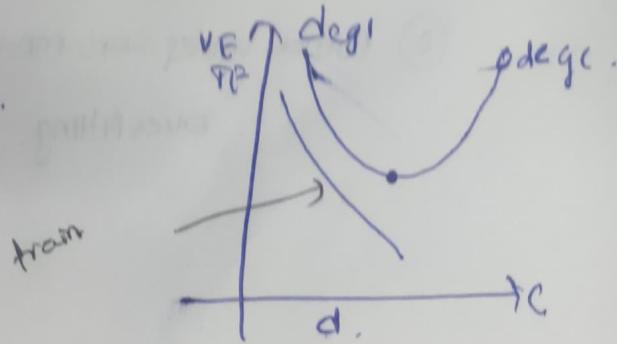
→ dcp (d)

→ low dp (eg  $d=1$ ) might underfit the data

→ high dp (eg  $d=10$ ) might overfit the data

→ Tuning helps find the optimal degree  
the balances bias & variance.

Degree: hyperparameters  
Coefficient parameter.



## Regularization parameter ( $\alpha$ )

- used in ridge (L2) or Lasso (L1) regression
- helps prevent overfitting by penalizing large coefficients
- higher  $\alpha$  - more regularization

## HPT methods

- ① Grid search
- ② Random search
- ③ cross-validation

⑧

## HPT def

### → key parameters

① max\_depth → limits tree depth to prevent overfitting

① low values → underfitting

② high values → overfitting

② min\_samples\_split → min samples required to split a node

① high values → prevents unnecessary splits

② lower values → more splits, leading to overfitting.

- ③ min\_samples\_leaf  $\rightarrow$  min samples per leaf node
- ① HV  $\rightarrow$  simpler model, reduces overfitting
  - ② LV  $\rightarrow$  more complex tree, risk of overfitting
- ④ max\_features  $\rightarrow$  limits no. of features considered per split
- Best practice:  $\rightarrow$  use  $\sqrt{n\text{-features}}$  for classification,  $\log_2(n\text{-features})$  for regression
- ⑤ criterion  $\rightarrow$  split quality metric
  - $\rightarrow$  gini or entropy for classification
  - $\rightarrow$  mse for regression
- ⑥ ccp\_alpha  $\rightarrow$  post-pruning parameter
  - $\rightarrow$  Higher alpha  $\rightarrow$  more pruning
  - $\rightarrow$  lower alpha  $\rightarrow$  less pruning

## HPT Methods

- ① Grid search
- ② Random Search
- ③ Bayesian Optimization
- ④ Cross-validation

## 5) Cross-validation

### Over-fitting

→ model learns both the patterns & noise from training data, reducing generalization ability

→ High variance, low bias

→ poor accuracy due to lack of generalization

→ too complex

(too many features)

### Underfitting

→ model is too simple to capture patterns in the data leading to poor performance

→ High bias, low variance

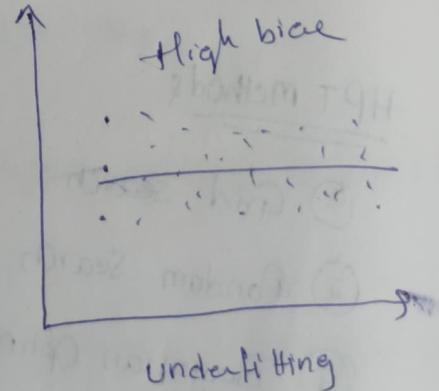
→ poor accuracy due to overly simplistic model

→ too simple

(insufficient features)

e.g. - A deep neural network memorizing training data instead of generalizing.

e.g. - A linear regression model failing to fit a non-linear dataset



How cross-validation helps detect overfitting & underfitting

- CV helps checks if a model is too-complex (overfit) or too simple (underfit) by testing it on different parts of the data

### Detecting issues

→ overfitting :- def

→ underfitting :- def

### Fixing

over

→ Use regularization

→ Reduce model complexity

→ Use more training data  
or data augmentation

under

→ Use a more complex model

→ Reduce regularization

→ Improve feature selection & preprocessing

### ④ K-fold CV

- It is a technique used to evaluate ML models by splitting the dataset into  $k$  equal parts (folds). The model is trained on  $k-1$  folds and tested on the remaining 1 fold, repeating this process  $k$  times with different test folds. The final performance is the average of all  $k$ -iterations.

## How it helps

- Reduce bias by ensuring every data point gets a chance in both test & training data sets.
- Helps detect overfitting & underfitting
- ensures that the model is tested on diff subsets of the data.

## k-fold cross-validation

- splits data into k-folds & runs k-train/testing iterations

- reduce variance by averaging results over k-runs

- uses more data for training, making the model better

- slower computation time

## Train-test split

- splits data into one training set & one test set  
e.g. - (80-20) split

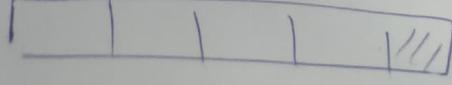
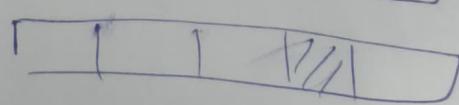
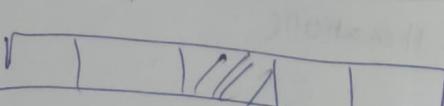
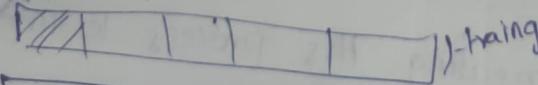
- can have high variance as the performance depends on a single test set

- less training data, as part of the dataset is reserved for testing

- faster computation time

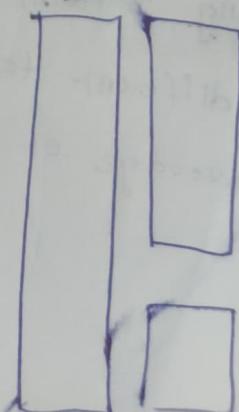
Dataset

100%



Dataset

-80%  
training



20%  
testing

## Why CV

- small datasets
- Hyperparameter tuning
- avoiding overfitting
- Helps detecting overfitting & underfitting
- reduces variance in performance metrics
- more reliable model evaluation.

## 16 Confusion matrix

→ It is a table to evaluate the performance of classification model by comparing actual vs predicted values

		Predicted +ve	-ve
Actual +ve	TP	FN	
-ve	FP	TN	

## Significance

- helps identify type of error
- essential for imbalanced datasets where accuracy alone is misleading.
- enables calculations of multiple performance metrics for better decision making.

TP - correctly predicted positive cases

TN - correctly -ve

FN - 2nc

FN - 2nc +ve (Type 1 error)

-ve (Type 2 error)

ii) precision, recall & f1-score.

① precision (positive predictive value)

↳ it measures how many of the predicted  
+ve cases were actually correct

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

High P  $\rightarrow$  fewer FP (useful in spam detection, where  
false alarms should be minimized)

low P  $\rightarrow$  Many incorrect +ve predictions

② Recall (sensitivity or TP Rate)

↳ it measures how many of the actual positive  
cases were correctly identified.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

High recall  $\rightarrow$  fewer false -ve (imp in medical diagnosis)  
Low  $\gamma$   $\rightarrow$  Many missed +ve cases

③  $f_1$ -score (Harmonic mean of precision & Recall)

It balances precision & recall, making it useful when both false +ve & false -ve matter

$$f_1\text{-score} = 2 \times \frac{P \times R}{P + R}$$

High  $f_1$   $\rightarrow$  good balance b/w P & R

Low  $f_1$   $\rightarrow$  indicates an imbalance b/w P & R

④ ROC (Receiver Operating characteristic)

$\rightarrow$  is a graphical tool used to evaluate the performance of a binary classifier at diff. threshold values.

$\rightarrow$  It helps assess how well a model distinguishes b/w +ve & -ve classes

$\rightarrow$  ROC curve plot

True +ve rate  $\Rightarrow$  TPR  $\Rightarrow$  Recall  $\rightarrow \frac{TP}{TP+FN}$

False -ve rate (FPR)  $\Rightarrow \frac{FP}{FP+TN}$

-) How to interpret the ROC-curve

④ Ideal curve

↳ The curve is close to the top-left corner.

↳ meaning high TPR & low FPR

⑤ Random model

↳ diagonal line ( $45^\circ$ ) represents random guessing ( $AUC = 0.5$ )

⑥ Poor model

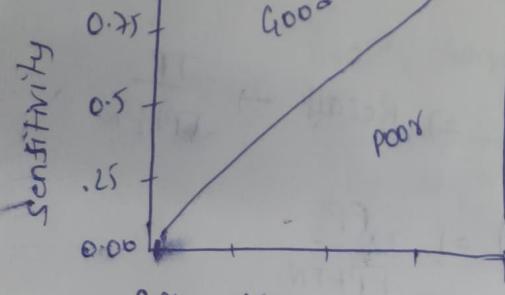
↳ below the diagonal means the classifier is worse than random.

-) AUC value model performance

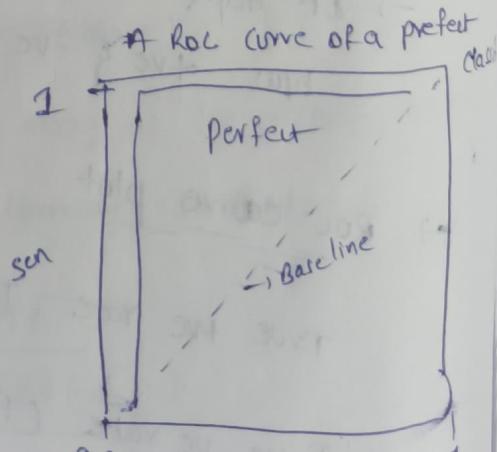
0.9 - 1 → Excellent

0.5 - 0.6 → Very poor.

A ROC curve of a random classifier



1 - specificity



1 - spe

## Use

- threshold selection
- imbalanced data handling
- comparing models

$$(14) \quad \underline{\text{diff}} \text{ b/w } \underline{R^2} \text{ & } \underline{\text{MSE}}$$

MSE (Mean squared error)

- measures the avg speed: squared diff b/w the actual  $y_i$  predicted values.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - \hat{y}_i)^2$$

lower MSE → better model fit

higher MSE → poor model performance

- used for comparing models in terms of absolute error

$R^2$  ( $R^2$ : co-efficient of determination)

- measures the proportion of variance in the dependent variable that is explained by the model

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

$SS_{res}$  → sum of squared errors

$SS_{tot}$  → total variance in the target variable

$R^2$  range 0 to 1

$R^2 = 1 \rightarrow$  perfect fit

$R^2 = 0 \rightarrow$  model is no better than the mean

$R^2 < 0 \rightarrow$  model is worse than using the mean as a predictor

-) Used for understanding how well the model explains variance in the data

MAE Mean absolute error (MAE) & MSE

MAE

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

→ measures the avg absolute diff. b/w actual & predicted values

⇒ lower MAE → better predictions

```
→ from sklearn.metrics import mean_absolute_error,  
    mean_squared_error
```

$$y_{\text{true}} = [3.05, 2.7]$$

$$y_{\text{pred}} = [2.5, 2.8]$$

```
mae = mean_absolute_error(y_true, y_pred)
```

```
mse = mean_squared_error(y_true, y_pred)
```

```
print(mae, mse)
```

→ without library

$$n = \text{len}(y_{\text{true}})$$

$$y_{\text{true}} = [ ]$$

$$y_{\text{pred}} = [ ]$$

```
mae = sum(abs(y_true[i] - y_pred[i]) for  
         i in range(n)) / n
```

```
mse = sum((y_true[i] - y_pred[i]) ** 2 for  
          i in range(n)) / n
```

```
print(mae)
```

```
print(mse)
```

18) Refer 13 Answer

$$13) \quad \text{accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

23 Why does linear SVM fail for non-linearly separable data.

a) Linear SVM tries to find a straight hyperplane that best separates two classes.

$\Rightarrow$  however, it fails when the data is non-linearly separable.

$\rightarrow$  No straight line can separate the classes

$\rightarrow$  if the data is non-linearly distributed, a linearly decision boundary cannot separate the two classes effectively

e.g. consider XOR data, where no straight line can separate the +ve & -ve classes

$\rightarrow$  High misclassification rate.

Since linear SVM forces a straight boundary, it will misclassify many points in a complex decision space

$\rightarrow$  poor generalization

i) It fails on unseen data when the actual distribution is non-linear

→ soln

- polynomial kernel
- radial basis function (RBF kernel)
  - ↳ maps data into infinite dimensions, making non-linear separation possible
- sigmoid kernel → mimics nn

23

## SVM

- ↳ aim to find a decision boundary (hyperplane) that maximally separates data points belonging to diff classes.

### Margin in SVM

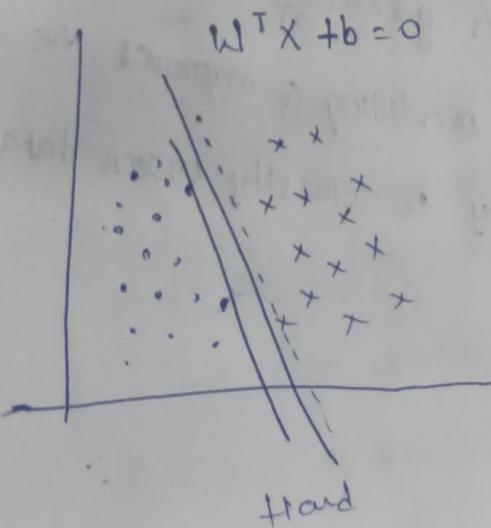
- ↳ Margin in SVM is the distance b/w the decision boundary (hyperplane) and the closest data points from each class

Large margin :- better generalization, a larger margin reduces overfitting & improves the model's ability to classify unseen data correctly.

Small margin  $\Rightarrow$  higher risk in overfitting, a small margin makes the classifier sensitive to noise

### types of margin.

	Hard margin	Soft Margin
applicability	used when data is perfectly separable	used when data is not perfectly separable
Misclassification	No	allows some
Margin	Maximized	balances margin size
overfitting risk	high	lower
use case	when training data is clean & linearly separable	When data has some noise or overlap b/w classes



## Optimization - frameworks for hard margin

Objective :- Maximize the margin

— the margin is given by  $\frac{1}{\|w\|}$ , so maximizing

the margin is equivalent to minimizing  $\|w\|$

→ given a dataset  $(x_i, y_i)$

where  $y_i \in \{-1, +1\}$

where  $w$  &  $b$  are weights & bias

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

### constraint

$$y_i(w \cdot x_i + b) \geq 1, \forall i$$

→ The constraint ensures that each data point lies outside or on the margin boundary

→ The Lagrange multipliers technique is used to solve this optimization problem.

## 24 diff kernels in SVM & their hyperparameters

→ SVM uses kernel functions to transform data into higher dimensions, allowing non-linearly separable data to be classified.

Kernel	Mathematical function	Hyperparameter	usage
linear kernel	$K(x, x') = x \cdot x'$	C (regularization parameter)	used when data is linearly separable
polynomial kernel	$K(x, x') = (x \cdot x' + b)^d$	C → degree of freedom b → coefficient	non-linear relationship with polynomial features
Radial basis function (RBF)	$K(x, x') = e^{-\gamma \ x - x'\ ^2}$	C $\gamma$ (gamma)	complex, non-linearly separable data
Sigmoid kernel	$K(x, x') = \tanh(\alpha x \cdot x' + c)$	C, C $\alpha$ (scaling factor)	similar to neural networks, used when data is binary
Gaussian kernel	$e^{-\frac{\ x - x'\ ^2}{2\sigma^2}}$	C $\sigma$ (sigma)	high-dimensional problems

## g) SVD (singular value decomposition)

↳ is a matrix factorization technique that decomposes any given matrix  $A$  into three matrices

$$A = U \Sigma V^T$$

where  
 $U \Rightarrow$  Orthogonal matrix containing the left singular vectors

$\Sigma \Rightarrow$  diagonal matrix with singular values

$V^T \Rightarrow$  orthogonal matrix [right] singular vector

→ used to in dimensionality reduction, noise reduction  
or data compression.

## How does SVD relate to PCA

→ PCA is a technique used for dimensionality reduction.  
it finds the principal components of the data  
which are the dirn of maximum variance.

① PCA is performed using SVD

→ PCA computes the eigenvectors of the co-variance matrix, but instead of calculating the co-variance matrix explicitly, it can

be obtained directly using SVD

→ the right singular vectors ( $V$ ) in SVD correspond to the principal components in PCA.

→  $\sigma_i$  represents the importance of each principal component

## ② dimensionality reduction

In PCA, we select the top  $k$  singular values & their corresponding singular vectors to reduce the no. of dimensions while preserving as much variance as possible.

$$A = \begin{bmatrix} 4 & 0 \\ 3 & -5 \end{bmatrix}$$

$$\left. \begin{array}{l} U \Rightarrow AA^T \\ V \Rightarrow A^TA \\ \sigma_i \Rightarrow A^TA \end{array} \right\} \begin{array}{l} \text{eigen vector} \\ \text{Eigen value} \\ \text{eigen vector} \end{array}$$

SVD :-

$$A = U \Sigma V^T$$

$S_2(U)$

$$AA^T = \begin{bmatrix} 4 & 3 \\ 0 & -5 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 3 & -5 \end{bmatrix}$$

$$= \begin{bmatrix} 25 & -15 \\ -15 & 25 \end{bmatrix}$$

Eigenvalues  $A^T A$

$$\det(A^T A - \lambda I) = 0$$

$$\begin{vmatrix} 25-\lambda & -15 \\ -15 & 25-\lambda \end{vmatrix} = 0$$

$$(25-\lambda)(25-\lambda) - (-15)(-15) = 0$$

$$\lambda^2 - 50\lambda + 400 = 0$$

$$\lambda_1 = 40, \lambda_2 = 10$$

$$\sigma_1 = \sqrt{40}$$

$$\sigma_2 = \sqrt{10}$$

$$Q = \begin{bmatrix} \sqrt{40} & 0 \\ 0 & \sqrt{10} \end{bmatrix}$$

$$A^T A = \begin{bmatrix} 25 & 15 \\ 15 & 25 \end{bmatrix}, \lambda_1 = 40$$

$$\begin{vmatrix} 25-40 & -15 \\ -15 & 25-40 \end{vmatrix} \begin{bmatrix} v_{11} \\ v_{12} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\lambda_2 = 10$$

$$v_2 = \frac{1}{\sqrt{2}} [1, 1]$$

$$v_1 = \frac{1}{\sqrt{2}} [1, -1]$$

$$V = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

Sy for U

$$AA^T = \begin{bmatrix} 40 \\ 35 \end{bmatrix} \begin{bmatrix} 4 & 3 \\ 0 & -5 \end{bmatrix} = \begin{bmatrix} 16 & 12 \\ 12 & 34 \end{bmatrix}$$

Eigen values

$$\lambda_1 = 40 \quad \lambda_2 = 10$$

$$\lambda_1 = 40$$

$$\begin{bmatrix} 16-40 & 12 \\ 12 & 34-40 \end{bmatrix} \begin{bmatrix} u_{11} \\ u_{12} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$u_1 = \frac{1}{\sqrt{13}} [4, 3]$$

$$\lambda_2 = 10$$

$$u_2 = \frac{1}{\sqrt{13}} [-3, 4]$$

$$U = \begin{bmatrix} \frac{4}{\sqrt{13}} & -\frac{3}{\sqrt{13}} \\ \frac{3}{\sqrt{13}} & \frac{4}{\sqrt{13}} \end{bmatrix}$$

Ortally

$$U = \begin{bmatrix} 4/\sqrt{10} & -3/\sqrt{10} \\ 3/\sqrt{10} & 4/\sqrt{10} \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 40 & 0 \\ 0 & 10 \end{bmatrix}$$

$$V = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

## 10) curse of dimensionality

↳ refers to the challenges & problems that arise when working with high-dimensional data.

→ Impact of increasing the no. of features

① increased sparsity

- ↳ In high-dimensional spaces, data points become more spread out, making it difficult for models to find meaningful patterns.

② Higher computational cost

- ↳ More features require more memory & processing power, slowing down training & inference.

③ Overfitting

- ↳ A larger no. of features can cause models to memorize noise instead of learning useful patterns, reducing generalization.

④ difficulty in dis-based algorithms

- ↳ Algo like K-NN & clustering rely on distance metrics, which become less effective in high dimensions as distances become almost uniform.

strategies to mitigate the curse of dimensionality

① feature selection

L, choose only the most relevant features

② dimensionality reduction

L, use techniques like PCA or SVD to transform high to low dimensional space while retaining imp info

③ Regularization

L, apply  $L_1$  or  $L_2$  regularization to prevent overfitting

④ collect more data

25

Eigen values

$$\lambda_1 = 4.5$$

$$\lambda_2 = 2.3$$

$$\lambda_3 = 0.7$$

Eigen vectors

$$v_1 = \begin{bmatrix} 0.8 \\ 0.5 \end{bmatrix} \quad v_2 = \begin{bmatrix} -0.4 \\ 0.9 \end{bmatrix} \quad v_3 = \begin{bmatrix} 0.2 \\ -0.1 \end{bmatrix}$$

dataset

$$x = \begin{bmatrix} 2 & 3 \\ 4 & 5 \\ 6 & 7 \end{bmatrix}$$

principal component

$$PC_1 = \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix} \quad PC_2 = \begin{bmatrix} -0.6 \\ 0.8 \end{bmatrix}$$

→ Identify the  $PC_1$  &  $PC_2$

$PC_1$  corresponds to the largest eigenvalue

$$\lambda_1 = 4.5$$

$$PC_1 = v_1 = \begin{bmatrix} 0.8 \\ 0.5 \end{bmatrix}$$

$$\lambda_2 = 2.3$$

$$PC_2 = v_2 = \begin{bmatrix} -0.4 \\ 0.9 \end{bmatrix}$$

compute projection on  $P_C$

↓ show step

projection =  $x \cdot P_C$

$$x = [x]$$

$$\begin{bmatrix} 2 & 3 \\ 4 & 5 \\ 6 & 7 \end{bmatrix} \begin{bmatrix} 0.3 \\ 0.6 \end{bmatrix} = \begin{bmatrix} 3.0 \\ 6.2 \\ 9.0 \end{bmatrix}$$

$\underline{P_C}$

$$\begin{bmatrix} 2 & 3 \\ 4 & 5 \\ 6 & 7 \end{bmatrix} \times \begin{bmatrix} 1.0 \\ 0.0 \end{bmatrix} = \begin{bmatrix} 2.0 \\ 4.0 \\ 6.0 \end{bmatrix} = x$$

$$\begin{bmatrix} 2 & 3 \\ 4 & 5 \\ 6 & 7 \end{bmatrix} \begin{bmatrix} -0.4 \\ 0.9 \end{bmatrix} = \begin{bmatrix} 1.2 \\ 1.5 \\ 2.0 \end{bmatrix}$$

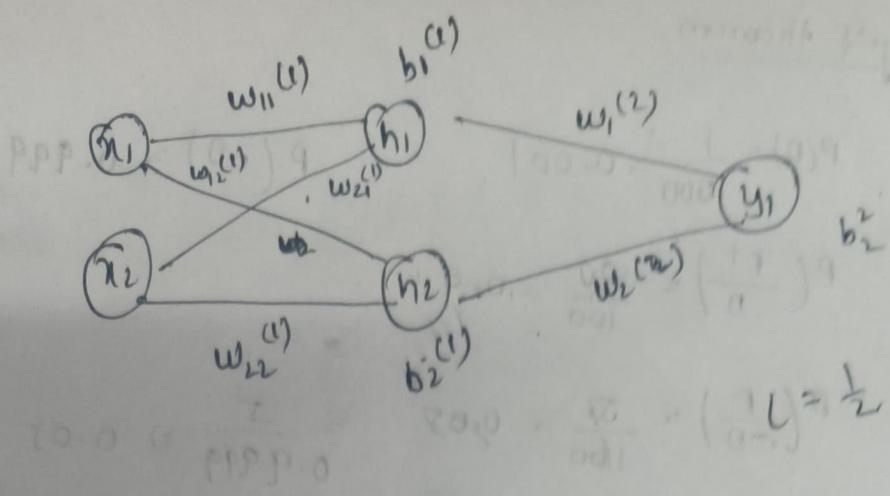
26

homework beginning

(1)  $x_1 = w_1^{(1)}$

(1)  $w_1^{(1)}$

(2)



$$z_1 = x_1(w_{11})^{(1)} + x_2(w_{21})^{(1)} + b_1^{(1)}$$

$$z_2 = x_2(w_{22})^{(1)} + x_1(w_{12})^{(1)} + b_2^{(1)}$$

$$h_1 = \sigma(z_1)$$

$$h_2 = \sigma(z_2)$$

$$z_3 = w_1^{(2)} h_1 + w_2^{(2)} h_2 + b_2^{(2)}$$

$$y = \sigma(z_3)$$

find  $w_{12}^{(1)}, w_{22}^{(1)}$

$$w_{12}^{(1)} \rightarrow z_2 \rightarrow h_2 \rightarrow z_3 \rightarrow y \rightarrow L$$

$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial z_3} \times \frac{\partial z_3}{\partial h_2} \times \frac{\partial h_2}{\partial z_2} \times \frac{\partial z_2}{\partial w_{12}}$$

$$w_{22}^{(1)} \rightarrow z_3 \rightarrow y \rightarrow L$$

$$\frac{\partial L}{\partial w_{22}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial z_3} \times \frac{\partial z_3}{\partial w_{22}}$$

$$\frac{\partial L}{\partial w_2} = 2(y - \hat{y}) \times \sigma'(z_3) \times \text{(check this)} w_2^{(2)}$$
$$\quad \quad \quad \times \sigma'(z_2) \times h_1$$

$$= 2(y - \hat{y}) \times \sigma'(z_3) \times h_2 + \cancel{w_2}$$

$$\frac{\partial L}{\partial w_1} = \cancel{y - \hat{y}}$$

## Bayes' theorem

$$P(D) = \frac{1}{1000} = 0.001$$

$$P(\neg D) = 0.999$$

$$P\left(\frac{T^+}{D}\right) = \frac{90}{100} = 0.9$$

$$P\left(\frac{F}{\neg D}\right) = \frac{5}{100} = 0.05 \quad = \frac{5}{0.999} = 0.005$$

$$P\left(\frac{D}{T}\right) = \frac{P(T^+/D) P(D)}{P(T^+/D) P(D) + P(F^-/D) P(\neg D)}$$

$$= \frac{(0.9)(0.001)}{(0.9)(0.0001) + (0.05)(0.999)}$$

$$= \frac{0.0009}{0.0001 + 0.049} = \frac{0.0009}{0.05009}$$

$$= \frac{0.0009}{0.0009} = \frac{0.0009}{0.05085} = 0.0177$$

$$P(\%T) = 1.77\%$$