

## Module 2 Regression (Ordinary least Squares (OLS)) / Polynomial regression / Ridge regression (regularization) / MLE / MAP

① training dataset  $D = \{(\vec{x}_i, y_i)\}_{i=1}^n$ ,  $\vec{x}_i \in \mathbb{R}^d$  feature vector input  
 $y_i \in \mathbb{R}$  target value; unknown mapping  $y = f(\vec{x})$ , hypothesis set output  
 $h_{\vec{w}}(\vec{x})$ , weight vector  $\vec{w} \in \mathbb{R}^d \rightarrow$  find  $\hat{y} = h_{\vec{w}}^*(\vec{x})$  best approx  $y = f(\vec{x})$  true value

②  $h_{\vec{w}}(\vec{x}) = w_0 + w_1 x_1 + \dots + w_d x_d$ ,  $\vec{x} = (x_1, \dots, x_d)^T$ ,  $\vec{w} = (w_0, w_1, \dots, w_d)^T$

$$h_{\vec{w}}(\vec{x}) = \vec{x}^T \vec{w} = \sum_{j=0}^d w_j x_j \quad \text{linear regression problem}$$

③  $\checkmark$  OLS:  $\hat{y}_i = h_{\vec{w}}(\vec{x}_i) = \vec{x}_i^T \vec{w}$ , we want  $\hat{y}_i \approx y_i$

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{pmatrix} = \begin{pmatrix} 1 & x_1 & \dots & x_d \\ 1 & x_1 & \dots & x_d \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \dots & x_d \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = X \vec{w}$$

1st data (1,  $x_1, \dots, x_d$ ) d num of feature

loss (cost) function (sum of squared error/SSE) target value in training data

$$L(\vec{w}) = \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \sum_{i=1}^n (h_{\vec{w}}(\vec{x}_i) - y_i)^2 = \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i)^2 = \|X \vec{w} - \vec{y}\|_2^2$$

optimization  $\vec{w}^* = \arg \min_{\vec{w}} \{L(\vec{w}) = \|X \vec{w} - \vec{y}\|_2^2\}$

$$L(\vec{w}) = \|X \vec{w} - \vec{y}\|_2^2 = (\vec{x} \vec{w} - \vec{y})^T (\vec{x} \vec{w} - \vec{y}) = \vec{w}^T X^T \vec{w} - 2 \vec{w}^T \vec{y} + \vec{y}^T \vec{y}$$

$$\nabla_{\vec{w}} L(\vec{w}) = 2X^T \vec{w} - 2\vec{y} = 0 \Rightarrow X^T \vec{w} - \vec{y} = 0 \Rightarrow \vec{w}^* = (X^T X)^{-1} X^T \vec{y}$$

$L$  is convex  $\leftarrow H(L(\vec{w})) = \sqrt{L(\vec{w})} = 2X^T X, \vec{w}^T (2X^T X) \vec{w} = 2(\vec{w}^T X \vec{w})^2 = 2\|\vec{w}\|_2^2 \geq 0$

④ Nonlinear Feature Mapping  $y = f(\vec{x})$  nonlinear

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d+1} (d>1), \vec{x}_i \in \mathbb{R}^d, \phi(\vec{x}_i) \in \mathbb{R}^{d+1}, \phi(\vec{x}_i)^T = (\phi_1(\vec{x}_i), \dots, \phi_m(\vec{x}_i))^T$$

hypothesis  $h_{\vec{w}}(\vec{x}) = \sum_{j=0}^d w_j \phi_j(\vec{x}) = \phi(\vec{x})^T \vec{w}$  basis function (polynomial, gaussian)

⑤  $\checkmark$  Polynomial Feature Mapping for 1 dimensional input

$$x \rightarrow \phi(x) = (1 x x^2 \dots x^m)^T, h_{\vec{w}}(x) = w_0 + w_1 x + \dots + w_m x^m = \phi(x)^T \vec{w}$$

input  $\rightarrow \vec{x}$

Dataset  $D = \{(\vec{x}_i, y_i)\}_{i=1}^n$ ,  $y = f(\vec{x})$ ,  $\hat{y} = h_{\vec{w}}(\vec{x}) = w_0 + w_1 x + \dots + w_m x^m = \phi(x)^T \vec{w}$

$$\vec{w} = (w_0, w_1, \dots, w_m)^T, \phi(x)^T = (1 x x^2 \dots x^m)$$

$$\begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \\ 1 & x_2 & x_2^2 & \dots & x_2^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^m \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{pmatrix} \quad \hat{y} = X \vec{w}$$

$$L(\vec{w}) = \|X \vec{w} - \vec{y}\|_2^2 \rightarrow \vec{w}^* = (X^T X)^{-1} X^T \vec{y}$$

⑥  $\checkmark$  Regularization - Ridge Regression ( $L_2$ ) (Lasso Regression ( $L_1$ )) validation

$$L(\vec{w}) = \|X \vec{w} - \vec{y}\|_2^2 + \lambda \|\vec{w}\|_2^2, \lambda > 0, \text{hyperparameter, choose through}$$

$$= \vec{w}^T X^T X \vec{w} - 2 \vec{w}^T \vec{y} + \vec{y}^T \vec{y} + \lambda \vec{w}^T \vec{w}$$

$$\nabla_{\vec{w}} L(\vec{w}) = 2X^T \vec{w} - 2\vec{y} + 2\lambda \vec{w} = 0 \Rightarrow \vec{w}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T \vec{y}$$

$$H(\vec{w}) = \vec{w}^T L(\vec{w}) = 2\vec{w}^T X + 2\lambda \vec{w} \text{ positive definite, minimum unique}$$

Regularization can be applied to overcome overfitting  $\rightarrow$  overfitting

⑦ Maximum Likelihood Estimate (MLE) Model

underlying model:  $y = f(\vec{x})$ , dataset  $D = \{(\vec{x}_i, y_i)\}_{i=1}^n$

$y_i \in \mathbb{R}$  is the observation of r.v.  $Y_i$ , i.e.  $Y_i = h_{\vec{w}}(\vec{x}_i) + \varepsilon_i$  — r.v.

assume  $\varepsilon_i \sim N(0, \sigma^2)$ , i.i.d. then  $Y_i \sim N(h_{\vec{w}}(\vec{x}_i), \sigma^2)$

$$L(\vec{w}, D) = P(Y_1=y_1, \dots, Y_n=y_n | \vec{x}_1, \dots, \vec{x}_n, \vec{w}) = \prod_{i=1}^n P(y_i | \vec{x}_i, \vec{w})$$

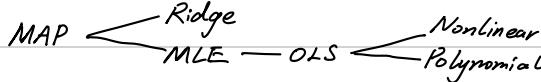
$$L(\vec{w}; X, y) = \ln(L(\vec{w}, D)) = \sum_{i=1}^n \ln(P(y_i | \vec{x}_i, \vec{w}))$$

$$P(Y_i=y_i | \vec{x}_i, \vec{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(y_i - h_{\vec{w}}(\vec{x}_i))^2}{\sigma^2}}$$

$$L(\vec{w}; X, \vec{y}) = -\sum_{i=1}^n \frac{(y_i - h_{\vec{w}}(\vec{x}_i))^2}{2\sigma^2} - N \ln(\sqrt{2\pi\sigma^2})$$

$$\begin{aligned} \vec{w}_{\text{MLE}} &= \arg \min_{\vec{w}} \left( \sum_{i=1}^n \frac{(y_i - h_{\vec{w}}(\vec{x}_i))^2}{2\sigma^2} + N \ln(\sqrt{2\pi\sigma^2}) \right) = \arg \min_{\vec{w}} \left( \sum_{i=1}^n (y_i - h_{\vec{w}}(\vec{x}_i))^2 \right) \\ (h_{\vec{w}}(\vec{x}) &= \vec{x}^T \vec{w}) = \arg \min_{\vec{w}} \{ \|X \vec{w} - \vec{y}\|_2^2 \} \quad (\text{OLS problem}) \end{aligned}$$

⑧ Maximum a Posterior (MAP) method Module 1 Pg 63



## Module 3 Performance Evaluation and Model Selection

(underfitting, overfitting / training, validation, test, cross-validation / performance metrics of regression / Bias - Variance trade-off)

① Underfitting: the model used is too simple to catch the trend behind the given dataset; overfitting: the model is so powerful it not only catch the trend but also the noise in the dataset.

② Example:  $\vec{w}_{\text{ridge}} = \arg \min_{\vec{w}} \{\|\vec{w}\|_2^2 + \lambda \|\vec{w}\|_1^2\} = (\vec{w}^T \vec{w} + \lambda I)^{-1} \vec{w}$

$\vec{w}$  — model parameter, a parameter to be optimized

order of polynomial model m, regularization weight  $\lambda$ , hyperparameters

③ true error  $R(\vec{w}) = \mathbb{E}_{\vec{x}, y} [L(h_{\vec{w}}(\vec{x}), y)]$ , expected loss over

the whole data distribution,  $\vec{w}^* = \arg \min_{\vec{w}} R(\vec{w})$  impossible

training error (in-sample error)  $\hat{R}_{\text{Train}}(\vec{w}) = \frac{1}{N} \sum_{i=1}^N L(h_{\vec{w}}(\vec{x}_i), y_i)$

test error (out-of-sample error)  $\hat{R}_{\text{test}}(\vec{w}) = \frac{1}{M} \sum_{i=1}^M L(h_{\vec{w}}(\vec{x}_i), y_i)$

④ The purpose of validation set is to determine the best value of the hyperparameters (used for hyperparameter tuning). Then, the tuned model will be tested using the test dataset.

⑤ Cross-validation: K-fold cross-validation: Pg 4

1. Shuffle the training data, and partition it into K blocks

2. Train the model on all the data except block i, use block i to evaluate the model. 3. Average the K validation error

⑥ Training Data — optimize the model parameter

Validation Data — optimize the model's hyperparameter

Test Data — evaluate the optimized model

⑦ Goodness-of-fit (performance metrics of regression model)

$$\text{Mean Squared Error (MSE)} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{Mean Absolute Error (MAE)} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$R^2 \text{ score} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad \text{most of time } [0, 1] \rightarrow 1 \text{ better}$$

(8) Bias - Variance Decomposition  $P_{27,31,32,38,39}$

$$\text{expected squared error } E(\vec{x}, h) = E[(h(\vec{x}, D) - Y)^2]$$

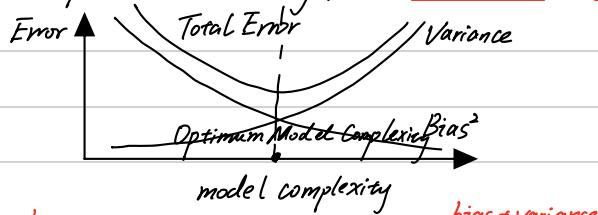
fixed test point  $\vec{x}$ , two r.v.  $D, Y$ , where  $Y = f(\vec{x}) + Z$

$$E(\vec{x}, h) = E[(h(\vec{x}, D) - Y)^2] = (E[h(\vec{x}, D)] - f(\vec{x}))^2 + \text{Bias} + \text{Variance} + \text{Noise}$$

Bias - Variance tradeoff (underfitted - Good fit/Robust - overfitted)

**Oversetting** - model captures the noise of the data, fits the data too well, shows low bias but high variance

**Underfitting** - model can not capture the underlying trend of the data does not fit the data well enough, shows low variance but high bias



Training error reflects bias not variance. Test error reflects both.

{ Decreasing the bias will increase the variance }

{ Decreasing the variance will increase the bias }

Use validation or cross-validation to tune some hyperparameters of your algorithm to achieve the trade-off

Regularization to overcome overfitting, bias ↑ variance ↓

Example Module 3 P42-49

(9) Reduce avoidable bias  $\leftarrow$  increase model size (complexity)  
Reduce/eliminate regularization

Reduce variance - Decrease the model size (complexity)/  
add regularization / add more training data / Feature selection  
to reduce number of input features / add early stop during training

Module 4 Perception and Logistic Regression (classification/performance metrics for classifiers / Perceptron model / Gradient Descent / Perceptron training algorithm (batch training / sequential training) / Logistic Regression )

① Classification Module 4 P1-6

Decision regions, Decision boundaries

Discriminant function 判别函数

② Performance Metrics for Classifier

Confusion Matrix  $\begin{cases} \text{multi-class classifier case} \\ \text{binary classifier case} \end{cases}$

P15		Actual class		TP - True Positive
Classifier		Positive	Negative	TN - True Negative
Predicted class	Positive	TP	FP	FP - False Positive
	Negative	FN	TN	FN - False Negative
		Accuracy = $\frac{TP+TN}{TP+FP+FN+TN}$ (imbalance dataset, mislead)		

$$\text{Error rate} = 1 - \text{Accuracy}$$

P15-24

$$\text{Precision} = \frac{TP}{TP+FP} \quad \text{Recall} = \frac{TP}{TP+FN}$$

$$F_1 = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{2 \cdot \text{recall} \cdot \text{precision}}{\text{precision} + \text{recall}} \quad F_B \text{ score}$$

③ linearly separable, line/hyperplane

④ Perception Model

Binary linearly separable case: input  $\vec{x}$ , two classes  $C_1, C_2$ .

$y=1$  for  $C_1$ ,  $y=-1$  for  $C_2$ . discriminant function  $g(\vec{x}) = \vec{w}^T \vec{x}$

$\vec{w}^T \vec{x}_i > 0$ ,  $\vec{x}_i$  in  $C_1$ ;  $\vec{w}^T \vec{x}_i \leq 0$ ,  $\vec{x}_i$  in  $C_2$

$\vec{w}^T \vec{x}$  is called linear score of the input  $\vec{x}$

1 计算时不带  $w_0$   $y(\vec{x}) = f(\vec{w}^T \vec{x})$

$x_1, w_1, w_0$   $x_2, w_2$  weighted sum linear score  $\rightarrow$  sign function  $\rightarrow$

$x_d, w_d$   $-1, \text{if } x < 0$   
 $0, \text{if } x = 0$   
 $1, \text{if } x > 0$

$$\Rightarrow \vec{w}^T \vec{x}_i > 0, \vec{x}_i \in C_1 (y_i = 1) \quad \text{normalized linear score}$$

$$\vec{w}^T \vec{x}_i \leq 0, \vec{x}_i \in C_2 (y_i = -1) \Rightarrow \vec{w}^T \vec{x}_i \geq 0, \text{any } \vec{x}_i \in D$$

✓ Perception criterion: the error (loss) function

$$E_p(\vec{w}) = - \sum_{i \in M} \vec{w}^T \vec{x}_i y_i$$

M denote the set of all misclassified instances

try to minimize the quantity  $-\vec{w}^T \vec{x}_i y_i$

Gradient Descent Procedure

$$\vec{w}^{k+1} = \vec{w}^k - \eta \nabla E_p(\vec{w}^k)$$

indexes the steps      stepsize      learning rate      gradient vector

$$\nabla E_p(\vec{w}) = - \sum_{i \in M} \vec{x}_i y_i \Rightarrow \vec{w}^{k+1} = \vec{w}^k + \eta \sum_{i \in M} \vec{x}_i y_i$$

✓ The Perceptron training Algorithm

Batch size (批量大小) - 每次梯度更新时使用样本的数量

epoch (训练周期) - one cycle through the entire training dataset

• Batch Gradient Descent 全批量梯度下降

批量大小等于整个训练集大小

• Mini-batch Gradient Descent 小批量梯度下降

批量大小介于 1 和训练集大小之间 (如 32, 64)

• Stochastic Gradient descent, SGD 随机梯度下降

批量大小为 1, 速度快, 方向不准, 收敛路径曲折 shuffling the training set each epoch

## ✓ Batch Training Algorithm 批训练 P<sub>1</sub>

Batch size =  $N$ , 每次迭代中，批训练使用整个训练集的数据来计算误差并更新权重。所有数据遍历后 misclassified 的有在  $M_k$  中， $\vec{w} \leftarrow \vec{w} + \eta \sum_{i \in M_k} \vec{x}_i y_i$

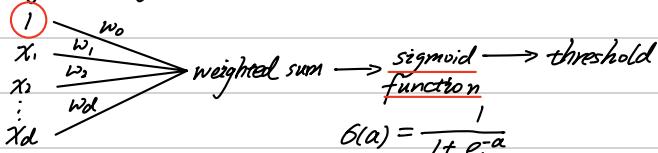
## ✓ Sequential Training Algorithm 顺序训练 P<sub>2</sub>

Batch size = 1, 每次迭代中使用单个样本来更新权重。  
 $K \leftarrow 0$ ,  $K \leftarrow K \bmod N$  (元余余数), if  $\vec{x}_k$  misclassified,  
 $\vec{w} \leftarrow \vec{w} + \eta \vec{x}_k y_k$

(5) logic OR, XOR problem P<sub>43, 44</sub>

(6) Sequential training Example P<sub>45-58</sub>

(7) Logistic Regression (Classification) P<sub>90-92</sub>



$$E(\vec{w}) = -\ln p(\vec{y}/\vec{w}) = -\sum_{i=1}^N \{y_i \ln p_i + (1-y_i)/n(1-p_i)\}$$

$$p_i = P(C_i|\vec{x}_i) = \sigma(\vec{w}^T \vec{x}_i) = \frac{1}{1 + e^{-\vec{w}^T \vec{x}_i}}$$

$$\nabla_{\vec{w}} E(\vec{w}) = -\sum_{i=1}^N \left( \frac{y_i}{p_i} - \frac{1-y_i}{1-p_i} \right) p_i(1-p_i) \vec{x}_i = -\sum_{i=1}^N (y_i - p_i) \vec{x}_i$$

$$\vec{w}_{k+1} = \vec{w}_k - \eta \nabla_{\vec{w}} E(\vec{w}) = \vec{w}_k + \eta \sum_{i=1}^N (y_i - p_i) \vec{x}_i$$

$$\text{where } p_i = \sigma(\vec{w}_k^T \vec{x}_i) = \frac{1}{1 + e^{-\vec{w}_k^T \vec{x}_i}}$$

• Regularized Logistic Regression  $E(\vec{w}) = -\ln p(\vec{y}/\vec{w}) + \lambda \|\vec{w}\|_2^2$  P<sub>84</sub>

• Multiclass Logistic Regression

## Module 5 SVM (margin/SVM/The lagrange multiplier method for constrained optimization/The lagrange dual problem)

(1) Margin of hyperplane — measure the distance between the hyperplane and the closest data points from each class.

hyperplane  $H: \{\vec{x}: \vec{w}^T \vec{x} + b = 0\}$ , find the distance from a point  $\vec{x}$  to hyperplane  $H$ .  $\vec{x} = \vec{x}_p + r \frac{\vec{w}}{\|\vec{w}\|}$

$r$  is the algebraic distance from  $\vec{x}$  to hyperplane  $H$

$$\vec{w}^T \vec{x} + b = \vec{w}^T \vec{x}_p + b + r \frac{\vec{w}^T \vec{w}}{\|\vec{w}\|} \Rightarrow r = \frac{\vec{w}^T \vec{x} + b}{\|\vec{w}\|}$$

$$\text{hyperplane } H: \{\vec{x}: \vec{w}^T \vec{x} + b = 0\}$$

algebraic distance of a point  $\vec{x}_i$  to  $H$  is  $r_i = \frac{\vec{w}^T \vec{x}_i + b}{\|\vec{w}\|}$

geometric distance of a point  $\vec{x}_i$  to  $H$  is  $\frac{|y_i(\vec{w}^T \vec{x}_i + b)|}{\|\vec{w}\|}$

The margin of hyperplane  $H$  is the geometric distance of the closest point in the data set to the hyperplane

$$\min_r \left\{ \frac{|y_i(\vec{w}^T \vec{x}_i + b)|}{\|\vec{w}\|} \right\}$$

Those point  $\vec{x}_i$  that satisfies  $\vec{w}^T \phi(\vec{x}_i) + b = \pm 1$  are called support vectors

The maximum margin solution is found by solving the optimization problem

$$\underset{\vec{w}, b}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\vec{w}\|^2 \right\}$$

subject to  $y_i (\vec{w}^T \phi(\vec{x}_i) + b) \geq 1, i = 1, 2, \dots, N$

## ✓ Lagrange Multiplier Method

$\min_{\vec{x} \in \mathbb{R}^n} f(\vec{x})$  subject to  $\begin{cases} C_i(\vec{x}) = 0, i \in \mathcal{E} \\ C_j(\vec{x}) \geq 0, j \in \mathcal{I} \end{cases}$

objective function      equality constraints      inequality constraint

Lagrangian function  $L(\vec{x}, \lambda) = f(\vec{x}) - \lambda_1 C_1(\vec{x})$

$$\nabla_{\vec{x}} L(\vec{x}, \lambda) = \nabla f(\vec{x}) - \lambda_1 \nabla C_1(\vec{x})$$

$$\nabla f(\vec{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(\vec{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\vec{x}) \end{pmatrix} \quad \nabla_{\vec{x}} L(\vec{x}, \lambda) = 0 \quad (\text{like } \nabla_{x_1} L(\vec{x}, \lambda) = 0, \nabla_{x_2} L(\vec{x}, \lambda) = 0)$$

$\Rightarrow x_1, x_2$  use constraint  $\Rightarrow \lambda \Rightarrow x$

Example P<sub>1</sub> P<sub>7-29</sub> P<sub>30-33</sub>

$$\min \{x_1 + x_2\}, \text{s.t. } 2 - x_1 - x_2 \geq 0, x_1 \geq 0$$

$$L(\vec{x}, \lambda) = f(\vec{x}) - \lambda_1 C_1(\vec{x}) - \lambda_2 C_2(\vec{x})$$

$$= x_1 + x_2 - \lambda_1 (2 - x_1 - x_2) - \lambda_2 x_2$$

$$\frac{\partial L(\vec{x}, \lambda)}{\partial x_1} = 2\lambda_1 x_1 + 1 = 0 \Rightarrow x_1 = -\frac{1}{2\lambda_1}$$

$$\frac{\partial L(\vec{x}, \lambda)}{\partial x_2} = 2\lambda_2 x_2 + 1 - \lambda_1 = 0 \Rightarrow x_2 = \frac{\lambda_1 - 1}{2\lambda_2}$$

$\lambda_1$  is strictly positive only when  $C_1(\vec{x})$  is active, that is

$$C_1(\vec{x}) = 2 - x_1 - x_2 = 0$$

$\lambda_2$  is strictly positive only when  $C_2(\vec{x})$  is active, that is

$$C_2(\vec{x}) = x_2 = 0$$

代回 constraints

$$\Rightarrow 2 - x_1^2 = 0 \Rightarrow x_1^2 = 2 \Rightarrow \frac{1}{4\lambda_1^2} = 2 \Rightarrow \lambda_1 = \frac{1}{2\sqrt{2}}$$

$$\frac{\lambda_2 - 1}{2\lambda_1} = 0 \Rightarrow \lambda_2 = 1$$

$$\Rightarrow \vec{\lambda} = \left( \frac{1}{2\sqrt{2}}, 1 \right)^T$$

$$\vec{x} = (x_1, x_2)^T = \left( -\frac{1}{2\lambda_1}, \frac{\lambda_2 - 1}{2\lambda_1} \right) = (-\sqrt{2}, 0)$$

## ✓ Lagrange Dual Problem

$$\min_{\vec{x} \in \mathbb{R}^n} f(\vec{x}) \text{ subject to } C_i(\vec{x}) \geq 0, i = 1, \dots, m$$

$$L(\vec{x}, \lambda) = f(\vec{x}) - \sum_{i=1}^m \lambda_i C_i(\vec{x})$$

$$\nabla f(\vec{x}) = \lambda_1 \nabla C_1(\vec{x}) + \lambda_2 \nabla C_2(\vec{x}) + \dots + \lambda_m \nabla C_m(\vec{x})$$

$$C_i(\vec{x}) \geq 0, i = 1, \dots, m$$

$$\lambda_i C_i(\vec{x}) = 0, i = 1, \dots, m$$

Dual objective function  $Q(\vec{\lambda}) = \inf_{\vec{x}} L(\vec{x}, \vec{\lambda})$

Dual problem  $\max_{\vec{x}} g(\vec{x})$  subject to  $\vec{x} \geq 0$

Example P37-39

$$\min_{x_1, x_2} \left( \frac{1}{2}(x_1^2 + x_2^2) \right) \text{ s.t. } x_1 - 1 \geq 0$$

$$L(x_1, x_2, \lambda_1) = \frac{1}{2}(x_1^2 + x_2^2) - \lambda_1(x_1 - 1)$$

$$\begin{aligned} \frac{\partial L}{\partial x_1} &= x_1 - \lambda_1 = 0 \Rightarrow x_1 = \lambda_1 \\ \frac{\partial L}{\partial x_2} &= x_2 = 0 \end{aligned}$$

Dual objective function

$$g(\lambda_1) = \frac{1}{2}(\lambda_1^2 + 0) - \lambda_1(1, -1) = -\frac{1}{2}\lambda_1^2 + \lambda_1$$

Dual problem

$$\max_{\lambda_1 \geq 0} \left( -\frac{1}{2}\lambda_1^2 + \lambda_1 \right)$$

$$\frac{d(-\frac{1}{2}\lambda_1^2 + \lambda_1)}{d\lambda_1} = 0 \Rightarrow \lambda_1 = 1$$

$$\Rightarrow \vec{x}^* = (x_1, x_2)^T = (\lambda_1, 0)^T = (1, 0)$$

④ SVM

The maximum margin solution is found by solving

$$\arg\min_{\vec{w}, b} \left\{ \frac{1}{2} \|\vec{w}\|^2 \right\} \text{ subject to } y_i(\vec{w}^T \phi(\vec{x}_i) + b) \geq 1 \quad i=1, \dots, N$$

hard margin P48-51

soft margin P52-59

Kernels and Kernel Trick P60-69

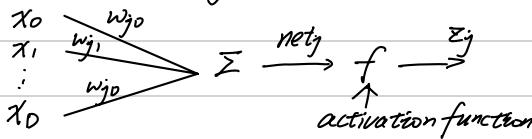
Module 6 ANN (Multi-layer Perceptron (MLP)) / forward propagation

backpropagation / regularization)

① Multiple-layer Perceptron Example P9-13

② ANN

1) The hidden layer

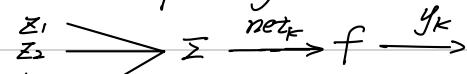


$$\text{net}_j = \sum_{i=0}^D x_i w_{ij} = \vec{w}_j^T \vec{x}$$

$$\vec{w}_j = (w_{j0} \dots w_{jD})^T, \vec{x} = (x_0, \dots, x_D)^T$$

$$z_j = f(\text{net}_j)$$

2) The output layer



$$\text{net}_k = \sum_{j=1}^M z_j w_{kj} = \vec{w}_k^T \vec{z}$$

$$y_k = f(\text{net}_k)$$

activation function P19

$$\text{sigmoid } f(\text{net}) = \frac{1}{1+e^{-\text{net}}}$$

$$\tanh f(\text{net}) = \frac{e^{\text{net}} - e^{-\text{net}}}{e^{\text{net}} + e^{-\text{net}}}$$

③ Training problem

$$E(\vec{w}) = \sum_{n=1}^N E_n(\vec{w}) = \frac{1}{2} \sum_{n=1}^N \|y(\vec{x}_n, \vec{w}) - \vec{t}_n\|^2$$

Gradient descent + backpropagation

$$\vec{w}^{t+1} = \vec{w}^t - \eta \nabla E(\vec{w}^t)$$

④ Backpropagation

At hidden layer, node j:

$$\text{net}_j = \sum_{i=0}^D x_i w_{ji} = \vec{w}_j^T \vec{x}, \quad z_j = f(\text{net}_j)$$

$$\begin{pmatrix} \text{net}_1 \\ \vdots \\ \text{net}_N \end{pmatrix} = \begin{pmatrix} w_{00} & w_{01} & \dots & w_{0D} \\ \vdots & \vdots & & \vdots \\ w_{N0} & w_{N1} & \dots & w_{ND} \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_D \end{pmatrix}$$

At output layer, node K

$$\text{net}_k = \sum_{j=1}^M z_j w_{kj} = \vec{w}_k^T \vec{z}, \quad y_k = f(\text{net}_k)$$

$$\begin{pmatrix} \text{net}_1 \\ \vdots \\ \text{net}_C \end{pmatrix} = \begin{pmatrix} w_{00} & \dots & w_{0m} \\ \vdots & \ddots & \vdots \\ w_{C0} & \dots & w_{Cm} \end{pmatrix} \begin{pmatrix} z_1 \\ \vdots \\ z_m \end{pmatrix}$$

instance  $(\vec{x}_n, \vec{t}_n)$  from training dataset, output  $\vec{y}_n$

$\vec{t}_n = \begin{pmatrix} t_1 \\ \vdots \\ t_C \end{pmatrix}$  target value ;  $\vec{y}_n = \begin{pmatrix} y_1 \\ \vdots \\ y_C \end{pmatrix}$  predicted value

Define the training error

$$E_n = \frac{1}{2} \sum_{k=1}^C (y_k - t_k)^2 = \frac{1}{2} \|\vec{y}_n - \vec{t}_n\|^2$$

First, hidden-to-output weights  $w_{kj}$

$$\frac{\partial E_n}{\partial w_{kj}} = \frac{\partial E_n}{\partial y_k} \cdot \frac{\partial y_k}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial w_{kj}} = (y_k - t_k) f'(\text{net}_k) z_j$$

sensitivity of node k,  $\delta_k = -\frac{\partial E_n}{\partial \text{net}_k} = (t_k - y_k) f'(\text{net}_k)$

$$\frac{\partial E_n}{\partial w_{kj}} = -\delta_k z_j$$

Second, input-to-hidden weights  $w_{ij}$

$$\frac{\partial E_n}{\partial w_{ij}} = \frac{\partial E_n}{\partial z_j} \cdot \frac{\partial z_j}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial w_{kj}} \cdot \frac{\partial w_{kj}}{\partial w_{ij}} = -\delta_j x_i$$

sensitivity for hidden node j  $\delta_j = \left( \sum_k w_{kj} \delta_k \right) f'(\text{net}_k)$

⑤ Improving Backpropagation training P36-42

1) Initializing weights → 2) Activation function → 3) Scaling input

4) Num of hidden nodes t) learning rate b) Overtraining problem

⑥ Regularization in Neural Network

$$\tilde{E}(\vec{w}) = E(\vec{w}) + \alpha \|\vec{w}\|^2$$

## Module 7 KNN and Naive Bayes Classifier (The k-Nearest Neighbor Method / Metrics (distance) / Curse of dimensionality / Naive Bayes Classifier)

### ① K-Nearest Neighbor (KNN) Module 7 P10 / Question P15

choose  $K$  - calculate distance - sort  $D$  - first  $k$ - $y$  model (if  $K=1$ )

### ② Metrics used in KNN: $\vec{x} = (x_1, \dots, x_m)^T$ , $\vec{y} = (y_1, \dots, y_m)^T$

Euclidean distance:  $d(\vec{x}, \vec{y}) = \left( \sum_{i=1}^m (x_i - y_i)^2 \right)^{\frac{1}{2}}$   $L_2$  norm

Minkowsky distance:  $d(\vec{x}, \vec{y}) = \left( \sum_{i=1}^m |x_i - y_i|^p \right)^{\frac{1}{p}}$   $L_p$  norm

Manhattan distance:  $d(\vec{x}, \vec{y}) = \sum_{i=1}^m |x_i - y_i|$   $L_1$  norm

Chebyshev distance:  $d(\vec{x}, \vec{y}) = \max_i |x_i - y_i|$   $L_\infty$  norm

Hamming distance:  $d_i = \begin{cases} 1, & \text{if } x_i \neq y_i \\ 0, & \text{otherwise} \end{cases}$ ;  $d(\vec{x}, \vec{y}) = \sum_i d_i$

cosine similarity measure  $\text{sim}(\vec{x}, \vec{y})$ ; correlation  $\text{corr}(\vec{x}, \vec{y})$  P14

### ③ Curse of dimensionality P17

$N$  data points,  $d$  dimensional unit hypercube, a hypercube of edge  $l$  contains  $K$  nearest neighbor of test set, fraction  $r = \frac{K}{N}$  ( $0 < l < 1$ )

$(d = r \Rightarrow l = r^{\frac{1}{d}})$ ,  $d \rightarrow \infty, l \rightarrow 1$  close to entire unit hypercube

数据稀疏性 Data Sparsity

This breaks down the KNN assumptions, because in high dimensions, the points in  $K$  nearest neighbor are not particularly closer than any other data points in the training set.

Dimension reduction (PCA) may provide a solution to alleviate the problem.

### ④ Naive Bayes Classifier - generative approach P5-29

training dataset  $D = \{(\vec{x}_i, y_i)\}_{i=1}^N$  belongs to class  $C_k$  if  $y_i = k$

new instance  $\vec{x} = (x_1, x_2, \dots, x_d)^T$ , predict the class of this instance

Bayesian approach to assign the most probable target value

$$\begin{aligned} y_{MAP} &= \arg \max_k P(C_k | x_1, \dots, x_d) \quad \text{attribute values} \\ &= \arg \max_k \frac{P(x_1, \dots, x_d | C_k) P(C_k)}{P(x_1, \dots, x_d)} \\ &= \arg \max_k P(x_1, \dots, x_d | C_k) P(C_k) \end{aligned}$$

Assume the attributes are conditionally independent  $\rightarrow$  Naive

$$P(x_1, x_2, \dots, x_d | C_k) = P(x_1 | C_k) \cdots P(x_d | C_k) = \prod_{j=1}^d P(x_j | C_k)$$

$$\Rightarrow y_{MAP} = \arg \max_k P(C_k) \prod_{j=1}^d P(x_j | C_k)$$

计算不同  $C_k$  的  $P(C_k) \prod_{j=1}^d P(x_j | C_k)$  选择最大值的类 (计算  $P(C_k)$  和  $\prod_{j=1}^d P(x_j | C_k)$ )

### ⑤ Laplace smoothing (additive smoothing) ( $P(\text{yellow}|t) = 0$ ) P30-33

$$\text{to calculate conditional probability } P(x_i | C_k) = \frac{N_{x_i, C_k} + 1}{N_{C_k} + d}$$

$N_{x_i, C_k}$  - 特征  $x_i$  在类别  $C_k$  训练集中出现的次数,  $d$  - dimensionality of the feature vector

$N_{C_k}$  - 训练集中属于类别  $C_k$  的样本数量

## Module 8 Decision Tree and Random Forest Model

(Entropy / Information Gain / Overfitting / Ensemble learning)

① Entropy - measure of impurity of a collection of samples (binary sets)  $\text{Entropy}(S) = -P_0 \log_2 P_0 - P_1 \log_2 P_1$  ( $0 \sim 1$ )

( $0 \log_2 0 = 0$ ) entropy  $\downarrow$  impurity  $\downarrow$ ,  $P_0 = 0.5$  entropy = 1  
(multi-class sets)  $\text{Entropy}(S) = -\sum_{i=1}^k P_i \log_2 P_i$

② Information gain (expected reduction in entropy)

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

feature 括号总数

③ Overfitting in Decision Tree: Limit the depth and nodes;

Post-prune the tree; use multiple trees to form a forest

④ Ensemble learning - Random Forest

Bagging (Bootstrap aggregating)

Random subset of features (attributes)

$L_2$  norm  $\vec{x} = (x_1, \dots, x_n)$ ,  $\|\vec{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$

$\underset{\vec{w}}{\text{argmin}} L(\vec{w})$  使  $L(\vec{w})$  取得最小值时  $\vec{w}$  的值

$$\nabla_{\vec{w}} (\vec{x}^T \vec{w}) = \vec{B}, \quad \nabla_{\vec{w}} (\vec{x}^T A \vec{w}) = 2A \vec{x}, \quad \nabla_{\vec{w}} (A \vec{w}) = A$$

$$\nabla f(x_1, \dots, x_n) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$