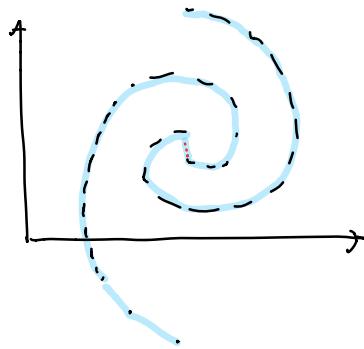


steps

single linkage example $x \in \mathbb{R}^2$



3.24 LEC
 k-means clustering
 avg/center
 ↑
 # of clusters

Input: n data points $x_1 \dots x_n$

$$\text{where } x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{id} \end{pmatrix} \in \mathbb{R}^d$$

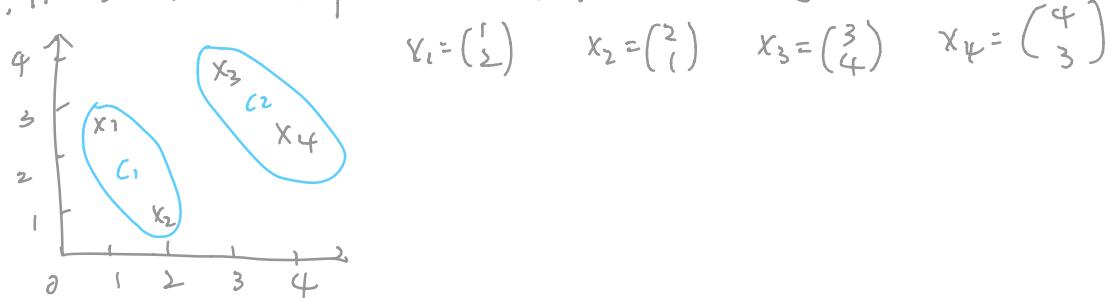
k-means Algorithm:

1. Init k cluster centers $c_1 \dots c_k \in \mathbb{R}^d$ (not overlapping)
2. Repeat until no changes in cluster assignment
3. For each item x_i , $i=1 \dots n$, find its closest cluster center. Let $y_i \in \{1, 2 \dots, k\}$ be that center
4. For each cluster center c_j , $j=1 \dots k$

$$c_j = \frac{\sum_{i=y_i=j} x_i}{\sum \dots 1}$$

$i=y_i=j$

Ex. in 2D with $n=4$ points and $k=2$ clusters



k-means

① step 1 arbitrary $c_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ $c_2 = \begin{pmatrix} 4 \\ 1 \end{pmatrix}$

Iteration 1:

$$d(x_i, c_1) = |x_i - c_1|$$

$$y_1=1 \quad (x_1 \text{ belongs to } c_1) \quad y_2=1 \quad y_3=1 \quad y_4=2$$

$$\text{Update } c_1 = \frac{\sum_{\substack{i: y_i=1 \\ i=1}} x_i}{1} = \frac{x_1 + x_2 + x_3}{3} = \left(\frac{\frac{1+2+3}{3}}{\frac{2+1+4}{3}} \right) = \begin{pmatrix} 2 \\ \frac{1}{3} \end{pmatrix}$$

$$c_2 = \begin{pmatrix} 4 \\ 1 \end{pmatrix}$$

Iteration 2:

$$y_1=1 \quad y_2=1 \quad y_3=2 \quad y_4=2$$

$$\text{Update } c_1 = \frac{x_1 + x_2}{2} = \begin{pmatrix} \frac{3}{2} \\ \frac{3}{2} \end{pmatrix} \quad c_2 = \frac{x_3 + x_4}{2} = \begin{pmatrix} \frac{7}{2} \\ \frac{1}{2} \end{pmatrix}$$

Iteration 3:

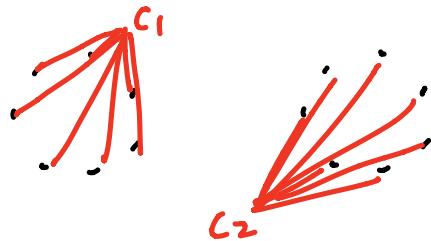
$$y_1=1 \quad y_2=1 \quad y_3=2 \quad y_4=2$$

Stop k-means, return $c_1 \dots c_k, y_1 \dots y_k$

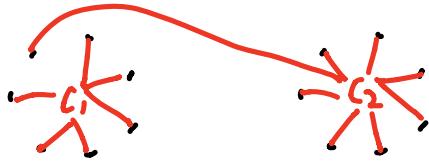
$$\sum_{i=1}^n \|x_i - c_{y_i}\|^2$$



good c_1, c_2
good $y_1 \dots y_n$



Bad c_1, c_2
good $y_1 \dots y_n$



good $c_1, c_2, y_2 \dots y_n$
bad y_1

The optimization problem of k-means

$$\min_{c_1 \dots c_k \in \mathbb{R}^d} \sum_{i=1}^n \|x_i - c_{y_i}\|^2$$

$$y_1 \dots y_n \in \{1 \dots k\}$$

Fix $c_1 \dots c_k$, solve $\min_{y_1 \dots y_n \in \{1 \dots k\}} \sum_{i=1}^n \|x_i - c_{y_i}\|^2$

$$\Leftrightarrow \min_{y_i \in \{1 \dots k\}} \|x_i - c_{y_i}\|^2 \quad \text{for } i=1 \dots n$$

\Leftrightarrow assign x_i to the closest cluster center
 $y_i = \arg \min_{j \in \{1 \dots k\}} \|x_i - c_j\| \quad \text{for } i=1 \dots n$

Fix $y_1 \dots y_n$, optimize $c_1 \dots c_k$

$$\min_{c_1 \dots c_k \in \mathbb{R}^d} \sum_{i=1}^n \|x_i - c_{y_i}\|^2 \Leftrightarrow \text{for } j=1 \dots k$$

$$\min_{c_j \in \mathbb{R}^d} \sum_{x_i = y_i = j} \|x_i - c_j\|^2$$

$$\begin{aligned}
 &\Leftrightarrow \nabla_{c_j} \sum_{x_i:y_i=j} \|x_i - c_j\|^2 \\
 &= \sum_{x_i:y_i=j} \nabla_{c_j} ((x_i - c_j)^T (x_i - c_j)) \\
 &= \sum_{x_i:y_i=j} 2(x_i - c_j)(-1) = \vec{0}
 \end{aligned}$$

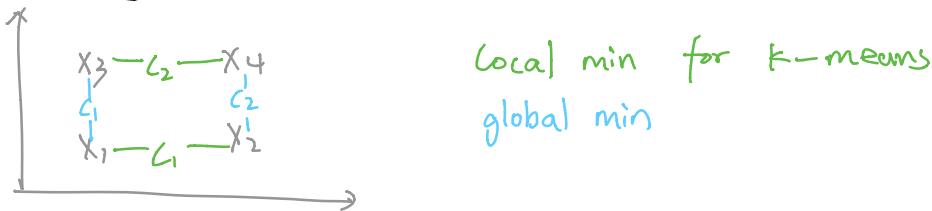
$$\sum_{x_i:y_i=j} (x_i - c_j) = \vec{0}$$

$$\sum_{x_i:y_i=j} x_i - \sum_{x_i:y_i=j} c_j = \vec{0}$$

$$\sum_{x_i:y_i=j} x_i = c_j \sum_{x_i:y_i=j} 1$$

$$c_j = \frac{\sum_{x_i:y_i=j} x_i}{\sum_{x_i:y_i=j} 1} \quad \text{avg of points assigned to } j$$

K-means Algorithm is a coordinate-descent alg to heuristically minimize the objective



Smart initialization of c_1, \dots, c_k : k-means++

- ① choose c_1 uniformly from x_1, \dots, x_n
- ② For $i=2 \dots k$

Let $D(x)$ be the distance from x to the closest initialized cluster center

D^2 weighting choose $c_j = x_i$ according to the distribution

$$p(x_i) = \frac{D^2(x_i)}{\sum_{i=1}^n D^2(x_i)} \quad i=1 \dots n$$

3.26 LEC Linear regression

1-D ordinary least squares (OLS) linear regression

Training set: $(x_1 \in \mathbb{R}, y_1 \in \mathbb{R}), \dots, (x_n, y_n)$

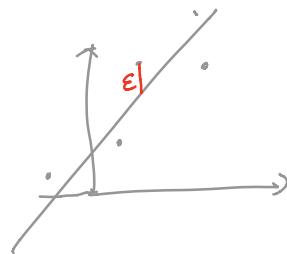
$\begin{matrix} \uparrow & \uparrow \\ \text{item} & \text{label} \\ \underbrace{\hspace{2cm}}_{\text{labeled item}} \end{matrix}$

Training: estimate a good model

Model: linear regression

$$y = \beta_0 + \beta_1 x + \varepsilon \quad \varepsilon \sim N(0, \sigma^2)$$

offset slope



Model parameter $\beta_0, \beta_1, \sigma^2$

Parameter estimation = training

view 1: maximum likelihood

view 2: least squares

Loss func: $l(x, y, \beta_0, \beta_1)$

least squares uses squared loss: $l(x, y, \beta_0, \beta_1) = (\beta_0 + \beta_1 x - y)^2$

Risk minimization

$$\text{Risk} = \sum_{i=1}^n l(x_i, y_i, \beta_0, \beta_1)$$

$$\hat{\beta}_0, \hat{\beta}_1 = \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^n \ell(x_i, y_i, \beta_0, \beta_1)$$

$$2^{\text{nd}} \text{ order: } y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$$

$$d^{\text{th}} \text{ order: } y = \beta_0 + \beta_1 x + \dots + \beta_d x^d + \epsilon$$

from input x , define a feature vector $\phi(x) = \begin{bmatrix} 1 \\ x \\ \vdots \\ x^d \end{bmatrix} \in \mathbb{R}^{d+1}$

$$\text{define } \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_d \end{bmatrix} \in \mathbb{R}^{d+1}$$

$$\text{Linear regression model } y = \phi(x)^T \beta + \epsilon$$

$$\text{Risk } \underbrace{\sum_{i=1}^n (\phi(x_i)^T \beta - y_i)^2}_{\text{risk}} \quad \text{or} \quad \frac{\sum_{i=1}^n (\phi(x_i)^T \beta - y_i)^2}{\text{MSE}}$$

Estimation (MLE)

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^{d+1}}{\operatorname{argmin}} \sum_{i=1}^n (\phi(x_i)^T \beta - y_i)^2$$

Define "design matrix" $X \uparrow_{\substack{\# \text{ items} \\ \# \text{ dim in } \phi(x)}}$

$$X = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix}_{n \times p} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}_{n \times 1}$$

$$n \times p \quad p \times 1$$

$$\text{model: } y = X\beta + \begin{pmatrix} \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\text{MLE: } \hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|X\beta - y\|^2$$

$$\nabla \|X\beta - y\|^2 = \vec{0}$$

$$2X^T(X\beta - y) = 0$$

$$X^T X \beta - X^T y = 0$$

$$\underbrace{X^T X}_{p \times p} \underbrace{\beta}_{p \times 1} = \underbrace{X^T y}_{p \times n}$$

assume $X^T X$ is invertible (full rank)

$$(X^T X)^{-1} (X^T X) \beta = (X^T X)^{-1} X^T y$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

Training set $(x_1, y_1), \dots, (x_n, y_n)$

Introduce test set $(x_{n+1}, y_{n+1}), \dots, (x_{n+m}, y_{n+m})$

(Both drawn from same distribution)

We train $\hat{\beta}$ on the training set

For each test item (x^*, y^*)

we predict its label as $\hat{y} = x^{*T} \hat{\beta}$

$$\text{we suffer loss } l(x^*, y^*, \hat{\beta}) = (\hat{y} - y^*)^2$$

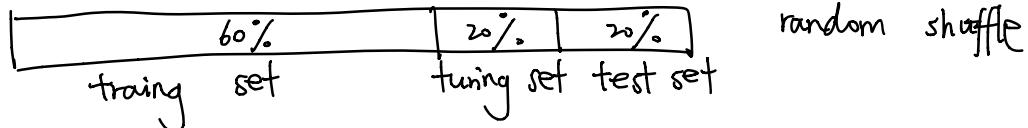
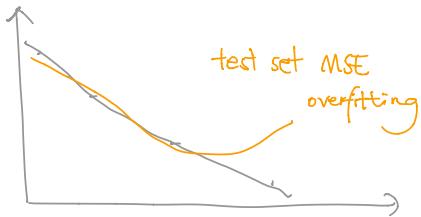
$$= (x^{*T} \hat{\beta} - y^*)^2$$

on the whole test set, we suffer

$$\text{MSE } \frac{1}{m} \sum_{i=n+1}^{n+m} l(x_i, y_i, \hat{\beta})$$

To minimize overfitting, use a tuning set

- ① get a labeled data set $(x_1, y_1) \dots (x_n, y_n)$
- ② randomly split the data set into 3 sets



- ③ Train $\hat{\beta}^{(0)}, \hat{\beta}^{(1)}, \dots, \hat{\beta}^{(n-1)}$ on training
- ④ Measure their tuning set MSE
pick the best model
 $j^* = \operatorname{argmin}_{j=0}^{n-1} \text{tuningset MSE}(\hat{\beta}^{(j)})$
we pick model $\hat{\beta}^{(j^*)}$
- ⑤ We report test-set MSE with model $\hat{\beta}^{(j^*)}$

3.3.1 LEC K nearest neighbor classifier

KNN algorithm:

- ↳ Input: $(x_1, y_1) \dots (x_n, y_n)$ where $x_i \in \mathbb{R}^d$, y_i a class label, a distance func $\text{dist}(x, x')$
- ↳ Given a new item $x \in \mathbb{R}^d$
- ↳ Find the k nearest neighbors of x in the training set under dist
- ↳ Predict a label \hat{y} as the majority label of the k nearest neighbors (break tie arbitrarily)

$$\text{lose}(x, u, \hat{u}) = 1 \text{ if } u \neq \hat{u} \quad \dots$$

$$\begin{aligned} & \text{--- } \cup \cup \cup \cup \quad \left\{ \begin{array}{ll} 0 & \text{if } \hat{y} = \check{y} \\ 1 & \text{if } \hat{y} \neq \check{y} \end{array} \right. \quad \text{0-1 loss} \\ & = \mathbb{1}[y \neq \hat{y}] \end{aligned}$$

Training set error (rate)

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}(x_i, y_i, \hat{y}_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[y_i \neq \hat{y}_i]$$

Test set error

$$\frac{1}{m} \sum_{i=n+1}^{n+m} \mathbb{1}[y_i \neq \hat{y}_i]$$

Why need a test set?

- Machine learning assumes an underlying joint distribution
- Training set is an independent and identically-distributed (i.i.d) sample from P

True error:

$$\mathbb{E}_{(x,y) \sim P(x,y)} \mathbb{1}[y \neq \hat{y}]$$

\hat{y} not known

$$\text{Accuracy} = 1 - \text{error}$$

How to choose k ?

Method 1: use training set

(x_1, y_1)	(x_n, y_n)	
training	tuning	test from training set

$$\hat{k} = \underset{k=1, 2, \dots}{\operatorname{argmin}} [\text{tuning error with respect to KNN predictions}]$$

report \hat{k}_{NN} prediction's test error

Method 2: cross validation

k -fold cross validation

For $i = 1 \dots k$

use fold i as the tuning set

all other sets as training set

get tuning error E_i

Pick model parameter (k in k_{NN})

$$\hat{k}_{k_{NN}} = \underset{k_{NN}}{\operatorname{argmin}} \frac{1}{k_{\text{fold}}} \sum_{i=1}^{k_{\text{fold}}} E_i$$

Retrain model on all folds as training set

Report test set error

4.2 LEC Regression

Logistic Regression

→ Input (training data)

$$(x_1, y_1) \dots (x_n, y_n)$$

$$\text{where } x_i = \begin{pmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{id} \end{pmatrix} \in \mathbb{R}^{d+1}$$

or $y_i \in \{-1, 1\}$ binary classification
 $y_i \in \{1, \dots, k\}$ k -classes

① Binary classification

$$P_w(y=1|x) = \frac{1}{1 + \exp(-x^T w)}$$

$$\text{if } x^T w = 0 \quad \frac{1}{1 + e^{-0}} = \frac{1}{2}$$

$$\text{if } x^T w = \infty \quad \frac{1}{1 + e^{-\infty}} = 1$$

