

Machine Learning

Lecture 1:

- * machine = computer
- * Learning = improving performance



⇒ A computer program improves its performance on a given task with experience

[1] Task:

- Classification** → map input into one of a set of classes [discrete]
- Regression** → map input into continuous output
- Transcription** → optical character recognition [OCR]
- Machine translation**

[2] Experience ⇒ what is the data (examples) that the program is using to improve performance?

[3] Performance measure ⇒ How is the performance evaluated?

a. **Accuracy**: produces the correct output [maximize]

b. **Error rate**: produces an incorrect output [minimize]

* **Supervised Learning**: $\begin{cases} \text{classification} \\ \text{regression} \end{cases}$

⇒ we are given a dataset and already know what our correct output should look like

* **unsupervised Learning**:

⇒ Little or no idea what our results should look like / no feedback

	Supervised	unsupervised
Discrete	classification	clustering
continuous	regression	Dimensionality reduction

* **Semi-supervised Learning**

↳ mix of supervised & unsupervised

يطلع الى امس معلوم من الجزء للعلوم

* **Reinforcement Learning**

↳ Model learns from a series of actions by maximizing a reward function

* Machine Learning depends on a number of algorithms:

$\begin{cases} \text{Computing resources available} \\ \text{nature of data} \end{cases}$

Lecture 2

* Linear regression

$$h(x) = \theta_0 + \theta_1 x \rightarrow \text{feature} \Rightarrow \text{hypothesis function}$$

y-intercept \leftarrow θ_0 \rightarrow slope \leftarrow θ_1

$$\frac{1}{2} \sum_{i=1}^m (y_i - h_{\theta}(x_i))^2 = \sum_{i=1}^m \hat{\epsilon}_i^2 \Rightarrow \text{Least squares}$$

actual \downarrow predicted \downarrow \rightarrow sum of the squared difference

كل الفرق يكون
المشرك على يبقى افضل

$$\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$$

$$= \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^i - y^i)^2$$

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^i - y^i)^2 \Rightarrow \text{cost function}$$

$$\text{at } \theta_1 = 2 = \frac{1}{2 \times 3} \begin{matrix} \theta_1 & x^i & y^i \\ (2 \times 1 - 1)^2 = 1^2 \\ (2 \times 2 - 2)^2 = 2^2 \\ (2 \times 3 - 3)^2 = 3^2 \end{matrix} = \frac{14}{6} = 2.33$$

no. of inputs \rightarrow

$$j=0: \frac{d}{d\theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)$$

$$j=1: \frac{d}{d\theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) \cdot x_i \Rightarrow \text{partial derivative}$$

repeat until convergence

$$\begin{cases} \text{new} & \text{old} \\ \theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \end{cases}$$

learning rate α

$$\theta_1 = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

\Rightarrow gradient decent

* provides continuous output

* How to find the best value of θ_1 ?

- \rightarrow plotting: unpractical specially in high dimensions
- \rightarrow analytical solution: not applicable for large datasets
- \rightarrow numerical solution ex: gradient descent
- \rightarrow underfitting gained

Lecture 3

* Linear regression with multiple values

$$\Rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

↳ the weights used by the model indicates the effect of each descriptive feature on the prediction returned by model.

⇒ Same gradient descent equations

① feature scaling: make sure features are on a similar scale

$$\rightarrow 0 \leq x_i \leq 1 \text{ [Mean normalization]} \quad x_i = \frac{x_i - \text{average}}{\text{Range} - \text{أقصى}} \quad \text{or } x_i = \frac{x_i - \text{أدنى}}{\text{أقصى} - \text{أدنى}}$$

② Learning rate [0.001 : 1]

α too small → slow convergence

α too large → may fail to converge = divergence

⇒ if gradient is working properly then $J(\theta)$ should decrease after each iteration

③ polynomial regression

* overfitting:

→ models the training data too well.

→ Learns the detail & noise in the training data to the extent that it negatively impacts the performance of the model on new data.

→ Ex: Decision tree

* underfitting:

→ poor performance

→ Ex: Linear regression

* Linear → minimize error

Logistics → maximize probability (Likelihood)

* gradient descent: $\theta_j = \theta_j + \alpha \sum_{i=1}^m (y_i - \frac{1}{1+e^{-\theta x_i}}) \cdot x_{ij}$

Lecture 4

* Logistic regression

p: probability of an event occurring

1-p: probability of an event not occurring

→ remove range restriction

$$\text{odds} = \frac{p_i}{1-p_i} \Rightarrow \log\left(\frac{p_i}{1-p_i}\right) = \sum_j x_{ij} \theta_j = x_i \theta$$

* after what we did in linear regression:

$$P(y_i | x_i, \theta) = \left(\frac{1}{1+e^{-\theta x_i}}\right)^{y_i} \left(1 - \frac{1}{1+e^{-\theta x_i}}\right)^{1-y_i}$$

↳ for all observations:

$$\frac{d}{d\theta} J(\theta) = \frac{d}{d\theta} \sum_{i=1}^m \left[y_i \log\left(\frac{1}{1+e^{-\theta x_i}}\right) + (1-y_i) \log\left(1 - \frac{1}{1+e^{-\theta x_i}}\right) \right]$$

$$= \sum_{i=1}^m \left(y_i - \frac{1}{1+e^{-\theta x_i}} \right) \cdot x_i$$

* provide discrete output

maximum likelihood estimation
MLE

$$\frac{p_i}{1-p_i} = e^{x_i \theta}$$

$$p_i = (1-p_i) e^{x_i \theta}$$

$$p_i = e^{x_i \theta} - p_i e^{x_i \theta}$$

$$p_i + p_i e^{x_i \theta} = e^{x_i \theta}$$

$$p_i (1 + e^{x_i \theta}) = e^{x_i \theta}$$

$$p_i = \frac{e^{x_i \theta}}{1 + e^{x_i \theta}} = \frac{1}{e^{-x_i \theta} + 1} = \frac{1}{e^{-x_i \theta} + 1}$$

Sigmoid [0,1]

Lecture 5

- * KNN \rightarrow smallest k
 - classification & regression algorithm
 - Lazy learner \rightarrow average of k
 - Supervised

\Rightarrow Eager learner

- * long time learning
less time classifying
- * Linear Regression
Decision tree

Lazy learner

- * less time learning
more time classifying
- * K. Nearest neighbors

* Euclidean distance = $\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$

\Rightarrow How to choose k?

- ① odd k value for the 2 classes
- ② k must not be a multiple of the number of classes
- ③ k too small \rightarrow overfitting [noise]
- ④ k too big \rightarrow mis-classify
- ⑤ Learn k by cross-validation

Advantages:

- ① Effective
- ② Non-linear
- ③ can add data seamlessly

Disadvantages:

- ① Sensitive to noise
- ② Large memory requirements

Lecture 6

- * Decision tree $\begin{cases} \text{supervised} \\ \text{classification} \end{cases}$

\rightarrow Entropy: $H(S) = -P_{(+)} \log_2 P_{(+)} - P_{(-)} \log_2 P_{(-)}$ [To measure impurity]

if equal 1 \rightarrow equal number of +ve & -ve examples

if equal 0 \rightarrow same class

Entropy \downarrow purity \uparrow

\rightarrow Information gain [tells us how important a given attribute of the feature vector is]

Information gain = entropy (parent) - [average entropy (children)]

$$\text{Gain}(S, A) = H(S) - \frac{\text{probability}}{\text{total}} * H(S_{\alpha_1}) - \frac{\text{probability}}{\text{total}} * H(S_{\alpha_2})$$

\rightarrow overfitting gained

① Stop growing the tree

② pre-pruning: stop the algorithm before it becomes a fully grown tree

③ post-pruning: trim the nodes of the decision tree / validation

Lecture 7

* Naive Bayes

- classification
- independent assumptions
- useful for very large datasets

* MAP rule [Discrete values]

$$\begin{matrix} \rightarrow P(\text{Yes} | x') \\ P(\text{No} | x') \end{matrix} \left. \vphantom{\begin{matrix} \rightarrow P(\text{Yes} | x') \\ P(\text{No} | x') \end{matrix}} \right\} \begin{matrix} \text{احسب} \\ \text{دول} \end{matrix} \Rightarrow \begin{matrix} \text{الاكبر هيكون} \\ \text{هو التي تير اختياره} \end{matrix}$$

* Continuous-valued features:

$$\begin{aligned} \hat{P}(x_j | c_i) &= \frac{1}{\sqrt{2\pi} \sigma_{ji}} \exp \left(- \frac{(x_j - \mu_{ji})^2}{2\sigma_{ji}^2} \right) \\ &= \frac{1}{\sqrt{2\pi} \times \text{variance}} \exp \left(- \frac{(x - \text{mean})^2}{2 \times \text{variance}} \right) \end{aligned}$$

* Zero conditional probability \rightarrow if no example contains the feature value

• m - estimate

$$\hat{P}(a_{jk} | c_i) = \frac{n_c + m_p}{n + m}$$

n_c : no. of samples of a specific feature

n : no. of samples of yes/no

P : values of specific feature

$m = 1$

* parametric methods summarize data into a fixed number of parameters whose count does not increase as the number of data points increase.

↳ Ex: Linear regression (single & multivariate)
Logistic regression
Naive Bayes

* Non parametric methods do not summarize data into a fixed number of parameters

↳ Ex: k-NN / Decision tree

* As model complexity increases, bias \downarrow & variance \uparrow

— maximum likelihood estimation

* ML: maximize the data likelihood given the model, $\arg \max_w P(\text{Data} | w)$

— maximum a posteriori

MAP: $\arg \max_w P(w | \text{Data})$

* for small training sets, naive Bayes generally is more accurate than Logistic regression

Lecture 8:

- SVM for linearly separable binary set
- used for classification and regression
- SVM decides that the best separating line is the line that bisects and is perpendicular to the connecting line.

* Margin in terms of w

$$w \cdot x_2 + b = 1$$

$$w \cdot x_1 + b = -1$$

$$w \cdot x_2 + b - (w \cdot x_1 + b) = 1 - (-1)$$

$$w \cdot x_2 + b - w \cdot x_1 - b = 2$$

$$w(x_2 - x_1) = 2 \text{ [not yet perpendicular]}$$

$$\rightarrow \frac{w}{|w|} (x_2 - x_1) = \frac{2}{|w|} \text{ [perpendicular]}$$

\hookrightarrow unit vector

* SVM as a minimization problem:

$$\bullet \text{ Maximizing } \frac{2}{|w|} = \text{minimizing } \frac{|w|}{2}$$

$$\bullet \text{ min: } \frac{1}{2} |w|^2 \Rightarrow \text{quadratic problem}$$

$$y_i (w x_i + b) \geq 1 \Rightarrow \text{Linear constrain}$$

$$\hookrightarrow \text{if } y = 1 \Rightarrow w x + b = 1$$

$$y = -1 \Rightarrow w x + b = -1$$

* Lagrange multiplier:

$$L_P \equiv \frac{1}{2} |w|^2 - \alpha [y_i (x_i \cdot w + b) - 1 \quad \forall i]$$

$$\equiv \frac{1}{2} |w|^2 - \sum_{i=1}^L \alpha_i [y_i (x_i \cdot w + b) - 1]$$

$$\equiv \frac{1}{2} |w|^2 - \sum_{i=1}^L \alpha_i y_i (x_i \cdot w + b) + \sum_{i=1}^L \alpha_i$$

$$\frac{\partial L_P}{\partial w} = 0 \Rightarrow w = \left[\sum_{i=1}^L \alpha_i y_i x_i \right]$$

$$\frac{\partial L_P}{\partial b} = 0 \Rightarrow \sum_{i=1}^L \alpha_i y_i = 0$$

Lecture 9

* properties of artificial neural network:

- Learning from experience
- Non-linearity
- Adaptivity
- Fault tolerance

* processing of ANN:

(1) Network topology:

- Single layer network: input layer is fully connected to the output layer [no hidden layer]
- Multi layer network: one or more layers between input and the output layer.

(2) Adjustments of weights or learning \Rightarrow modifying weights

(3) Activation functions

• Binary: $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$

• ReLU: $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$

• Sigmoid: $f(x) = \frac{1}{1 + e^{-x}} \quad [0, 1]$

• tanh: $f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1 \quad [-1, 1]$

• Softmax: $f(x) = \frac{e^{x_i}}{\sum_{j=1}^J e^{x_j}} \quad \text{for } i = 1, \dots, J$

* weight adaptation:

[1] $\omega(u+1) = \omega(u) + \eta [d(u) - y(u)] x(u)$
old weight derived inputs
predicted

[2] prediction error \rightarrow prediction output \rightarrow SOP \rightarrow input weights

$$E = \frac{1}{2} (d - y)^2$$

$$f(s) = \frac{1}{1 + e^{-s}}$$

$$S = \sum_j x_j w_j + b_i \quad w_1, w_2$$

\leadsto chain rule: $\frac{\partial E}{\partial y} \times \frac{\partial y}{\partial s} \times \frac{\partial s}{\partial w_i} = \frac{\partial E}{\partial w_i}$

① $\frac{\partial E}{\partial y} = y - d$

② $\frac{\partial y}{\partial s} = \frac{1}{1 + e^{-s}} \left(1 - \frac{1}{1 + e^{-s}}\right)$

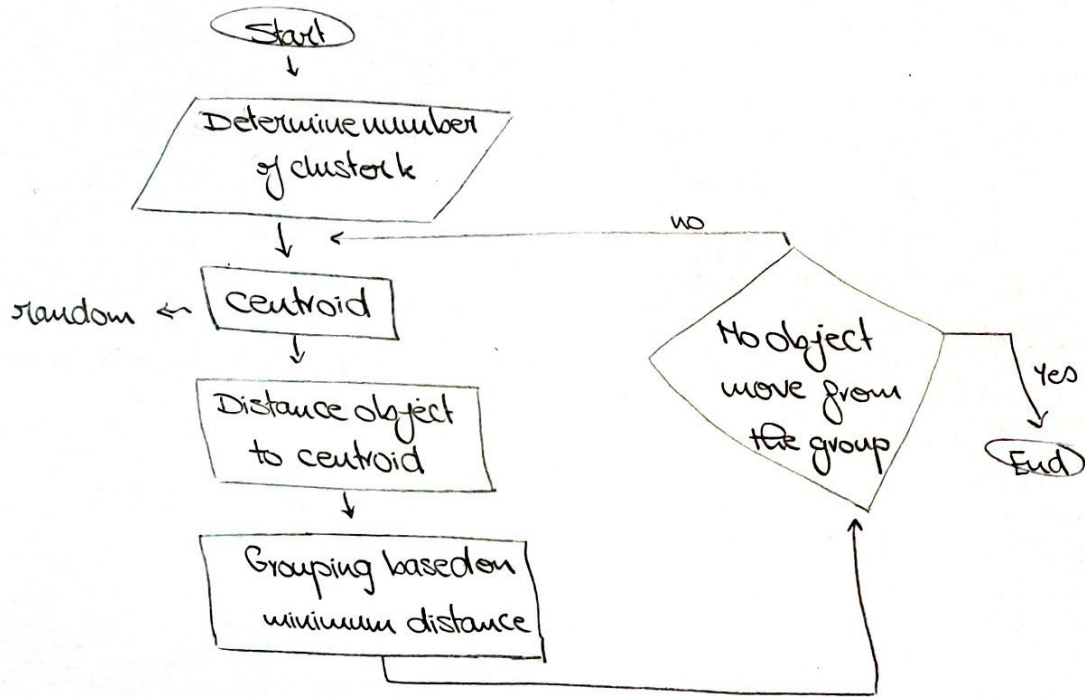
③ $\frac{\partial E}{\partial w_1} = x_1$ and $\frac{\partial E}{\partial w_2} = x_2$

④ $\frac{\partial E}{\partial w_i} = (y - d) \frac{1}{1 + e^{-s}} \left(1 - \frac{1}{1 + e^{-s}}\right) x_i$

$\hookrightarrow \omega_{i, \text{new}} = \omega_{i, \text{old}} - \eta * \frac{\partial E}{\partial w_i}$

Lecture 10

* cluster analysis is like classification, but the class label of each object is not known



* complexity: $O(n * k * I * d)$ \rightarrow no. of attributes
no. of points \downarrow \rightarrow no. of iterations
no. of clusters

* complexity of each round: $O(n * k)$

* Strength

\rightarrow relatively efficient

* Weakness

\rightarrow applicable only with numerical values

\rightarrow need to specify k

\rightarrow Sensitive to noise data

Lecture 11

* Two main types of hierarchical clustering:

- Agglomerative - merge
↳ Bottom-up approach [from 1 \rightarrow N]
- Divisive - split
↳ Top-down approach [from N \rightarrow 1]

* Distance calculations:

- Single-linkage clustering \rightarrow min
- Complete-linkage clustering \rightarrow max
- Average linkage
- centroid [distance between two centers]

* Complexity:

- Space: $O(n^2)$
- time: $O(n^3)$ or $O(n^2 \log(n))$

* Strength:

- no need to specify the number of clusters
- Easy to implement
- Dendrogram is very useful

* Weakness:

- can never undo any previous steps
- Time complexity is large
- Difficult with large dataset

Lecture 12

* Less number of features:

interpretability \uparrow accuracy \downarrow

* More number of features:

interpretability \downarrow accuracy \uparrow

* performance of ML model:

- ① choice of algorithm
- ② feature selection
- ③ feature creation
- ④ Model selection

* Feature selection methods:

① Filter:

- calculate the feature relevance score and remove low-scoring features
- Scalable, simple and fast. \checkmark
- each feature is considered separately \times

② Wrapper:

- Generate and evaluate various subsets of features. Based on accuracy
- feature dependency \checkmark
- very computationally intensive
- Ex: forward selection / backward elimination

③ Embedded:

- while the model is being created, it finds which is the best contribute to the accuracy
- Lasso: absolute value
- Ridge, Square

* precision: how many selected items relevant

Recall: how many relevant items selected?

$$F\text{-measure} = 2 * \frac{(\text{Precision} * \text{Recall})}{(\text{precision} + \text{Recall})}$$