

STK-INF4000 Notes - ML Basics

Dirk Hesse

Jan. 26th 2017

Intro

Basic idea: Some *data* collected and want to draw some *conclusions* about them.
E.g.

- Spam classification.
- Search result ranking.
- Sentiment analysis.
- Predictive maintenance.

Machine learning problems fall broadly into one of two categories:

- Supervised Learning
- Unsupervised Learning

Supervised Learning

- Random variable $X = (X_1, \dots, X_p)^T$.
- Output variable Y (usually *not* a vector, but could be) *or* G (group, e.g. $0, 1$).
- *Training examples* $(x_i, y_i), i = 1, \dots, n$.
- *Task*: Given a *value* for X , make a *good* prediction for Y , denoted usually \hat{Y} .

Examples

- Sentiment Analysis
- Demand Prediction
- Traffic Forecasting
- Drought Forecasting.
- Increasing Farm Yields.

What Can be Learned?

- Y *must* be dependent on X .
- Dependence can be *very* complex.
 - Deep learning.

Unsupervised Learning

No Y given. Just interested in properties of X .

Examples

- Grouping Things
- Recommendation
- Outlier Detection

Basic Probability

Random Variables

- Commonly denoted X or Y .
- On continuous or discrete spaces.
- Formally
 - $X : \Omega \rightarrow \mathbb{R}$
- PMF

$$f_X(x) = \mathbb{P}[X = x] = \mathbb{P}[\{\omega \in \Omega : X(\omega) = x\}]$$

- PDF

$$\mathbb{P}[a \leq X \leq b] = \int_a^b f(x)dx$$

- CDF

$$F_X : \mathbb{R} \rightarrow [0, 1], \quad F_X(x) = \mathbb{P}[X \leq x]$$

- Conditional:

$$f_{Y|X}(y|x) = \frac{f(x, y)}{f_X(x)}$$

Expectation Values

Discrete

$$\mathbb{E}[X] = \sum_x xp_X(x)$$

$$E[f(X)] = \sum_x f(x)p_X(x)$$

Continuous

$$E[X] = \int xp_X(x)dx$$

Conditional Probability

Study, 60% women, 40% men. High, low income.

	Female	Male	
High	9%	11%	20%
Low	46%	34%	80%
	55%	45%	

0.09, 0.11, 0.46, and 0.34 are called *joint* probabilities, while 0.55, 0.45, 0.2, and 0.8 are called *marginal* probabilities. The *conditional* probability $P(High|Female)$ can be calculated using the formula

$$P(x|y) = \frac{P(x,y)}{P(y)}$$

Here, $P(High|Female) \approx 16\%$, while $P(High|Male) \approx 24\%$.

Decision Theory - Supervised Learning

We have

- $X \in \mathbb{R}^p$ random input vector.
- $Y \in \mathbb{R}$ output vector.
- Looking for $f(X)$ predicting Y .

Why?

- Prediction
 - Exact form of f not too important.
 - Accuracy important.
 - * How good can we do?

- * Assume $Y = f(X) + \epsilon$.
- * Our estimate $\hat{f}(X) = \hat{Y}$:

$$E(Y - \hat{Y})^2 = E[f(X) + \epsilon - \hat{f}(X)]^2 = E[f(X) - \hat{f}(X)]^2 + \text{Var}(\epsilon).$$

- Inference
 - Which parts of X are important for predicting Y ?
 - What is the relationship between Y and X ?
 - What is the *functional* relationship between them? Linear?
 - **Exact form** of f is important.

How to get f ?

We need a *loss function* $L(Y, f(X))$, most commonly squared error loss, $L(Y, f(X)) = (Y - f(X))^2$.

How to choose f ?

$$\text{EPE}(f) = E(Y - f(x))^2 = \int (y - f(x))^2 \Pr(dx, dy)$$

Now $\Pr(dx, dy) = \Pr(Y|X) \Pr(X)$, so

$$\text{EPE}(f) = E_X E_{X|Y}([Y - f(X)]^2 | X)$$

and

$$f(x) = \underset{c}{\operatorname{argmin}} E_{Y|X}([Y - c]^2 | X = x)$$

This yields

$$f(x) = E(Y|X = x)$$

{#eq:fEYX}

K-Nearest Neighbors

KNN implements ({eq:fEYX}) in a very simple way:

$$\hat{f}(x) = \frac{1}{k} \sum_{z \in N_k(x)} z,$$

where $N_k(x)$ are the k closest training examples to x from a given training set.

Appendix

Formal Definitions

- Sample space Ω
- Outcome $\omega \in \Omega$
- Event $A \subseteq \Omega$
- σ -Algebra \mathcal{A}
- Probability distribution P
 - $P[A] \geq 0 \quad \forall A$
 - $P[\Omega] = 1$
 - $P[\bigcup_i A_i] = \sum_i P[A_i]$
- Conditional probability $P[A|B] = \frac{P[A \cap B]}{P[B]}$, if $P[B] > 0$.