



 Non-Parametric Classifiers, K-Nearest Neighbors - KNN
 Statistical Evaluation of Binary Classification: Confusion Matrix, ROC, AUC
 Parametric Probabilistic Classification: MLE & MAP Estimation, Bayes & Naive Bayes Classifiers

Prof. Vasilis Maglaris

<u>maglaris@netmode.ntua.gr</u>

www.netmode.ntua.gr

Room 002, New ECE Building

Tuesday May 13, 2025

NTUA - National Technical University of Athens, DSML - Data Science & Machine Learning Graduate Program

STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Classification based on *K*-Nearest Neighbors (1/2)

https://machinelearningmastery.com/parametric-and-nonparametric-machine-learning-algorithms/

- Non-Parametric Methods of Supervised Learning: Based on *distance* estimates from known sample vectors (*patterns*), without assuming a *probabilistic model* for the (labeled) training dataset that would determine the *form* of the function $y = h(\mathbf{x})$.
- By contrast, Parametric Methods are based on pre-selecting a *form* for the input-output function and learn its *coefficients* from the (labeled) training data (e.g. *Linear* & *Logistic Regression*, *Perceptron*, *Simple Neural Networks*)

K-Nearest Neighbors Algorithm (KNN)

Assume a *training* set { $\mathbf{x}(n), d(n)$ }, n = 1, 2, ..., N of *vectors* $\mathbf{x}(n) = [x_1(n) x_2(n) ... x_m(n)]^T$ with *labels* $d(n) \rightarrow C$ indicating the class of $\mathbf{x}(n)$ (e.g. $C \in \{0,1\}$ for binary classification) and a

distance *metric*, e.g. *Euclidean Distance* $||\mathbf{x}(i), \mathbf{x}(j)|| = \sqrt{\sum_{k=1}^{m} (x_k(i) - x_k(j))^2}$ or *Hamming Distance* - number of opposite binary digits between $(\mathbf{x}(i), \mathbf{x}(j))$

- Training Phase (Lazy Learning Method): Storage of training sample (N labeled patterns)
- **Test Phase**: A new vector $\hat{\mathbf{x}}$ is classified according to the *labels* of the *majority* of *K* nearest neighbors amongst the *N* training *patterns* (vectors) $\mathbf{x}(n)$

STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Classification based on *K*-Nearest Neighbors (2/2)

https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm

- *K*: Positive Integer **hyperparameter**, *odd* for binary classification, usually selected after *cross validation* trials
- K = 1: A new vector $\hat{\mathbf{x}}$ is classified according to the class of its nearest neighbor (*pattern*)
- $K \gg 1$: Tolerant in distortion/noise of classification regions *but* with increased storage requirements and prone to undesirable impact of *outliers*, i.e. rare patterns with extreme characteristics: Such cases require vector normalization, reduction of coordinates and filtering of *outliers*





Statistical Evaluation of Binary Classification: Confusion Matrix, ROC, AUC (1/3)

Statistical Binary Classification – Parametric Classifiers

Classification of examples (vectors, patterns) of a sample in 2 classes $C \in \{P, N\}$: **Positive P**, **Negative N**

- Diagnosis of infections: **Positive** \triangleq **Infected Sample Element**
- Identification of anomalies: Positive
 Anomalous Sample Element (outlier)
- Recognition of binary patterns (e.g. cats dogs): **Positive** \triangleq **Cat**, **Negative** \triangleq **Dog**
- Target identification signalling: **Positive** \triangleq **Foe**, **Negative** \triangleq **Friend** (**IFF**: Identify Friend or Foe)
- The **parametric** classification algorithms identify parameters of pre-selected input-output function $y = h(\mathbf{x})$, conforming to *training data*:
- Assumption of *Gaussian* distribution for the sample space to be identified with 2 parameters mean value - variance
- Examples of Parametric Methods:
 - > Parameter tuning in *logistic regression* from a *labeled* training sample
 - Parameter tuning of synaptic weights of *perceptron* neural networks via iterative *back-propagation* of output sample elements and convergence to statistically inferred properties from training datasets

https://machinelearningmastery.com/parametric-and-nonparametric-machine-learning-algorithms/ Parametric Methods: Linear & Logistic Regression, Perceptron Convergence, Bayes Classifiers... Non-Parametric Methods: KNN, Decision Trees, SVM...

Statistical Evaluation of Binary Classification: Confusion Matrix, ROC, AUC (2/3)

Confusion Matrix

- Incorrect Predictions: False Positives FP, False Negatives FN
- Correct Predictions: True Positives TP, True Negatives TN

Rates of Correct/Incorrect Predictions:

$$TPR = \frac{TP}{TP + FN}$$
, $TNR = \frac{TN}{TN + FP}$

$$FNR = \frac{FN}{FN + TP} = 1 - TPR$$
, $FPR = \frac{FP}{FP + TN} = 1 - TNR$



Confusion Matrix

Binary Classifier Evaluation Metrics

Accuracy: $ACC = \frac{TP+TN}{TP+TN+FP+FN}$ (ratio of total *correct* predictions) Sensitivity, Recall: $TPR = \frac{TP}{TP+FN}$ (*correct* positive predictions out of *actual* positive examples) Precision: $PRE = \frac{TP}{TP+FP}$ (*correct* positive predictions out of *positive* predictions) F1-Score: $F1 = \frac{TP}{TP+\frac{1}{2}(FP+FN)} = \frac{2}{TPR^{-1}+ACC^{-1}}$ (harmonic mean of *recall & precision*)

Statistical Evaluation of Binary Classification: Confusion Matrix, ROC, AUC (3/3)

Example of Image Recognition: Cat or Dog

https://en.wikipedia.org/wiki/Confusion matrix

Test Sample 12 images: 8 cats (class P), 4 dogs (class N)

The binary classifier predicts after training 7 cats and 5 dogs (9 correct, 3 incorrect) as shown in the *Confusion Matrix*:

- Accuracy: $ACC = \frac{TP+TN}{TP+TN+FP+FN} = \frac{6+3}{12} = 3/4$
- Sensitivity (Recall): $\mathbf{TPR} = \frac{\mathrm{TP}}{\mathrm{TP} + \mathrm{FN}} = \frac{6}{6+2} = 3/4$

Predicted class Actual class	Cat	Dog
Cat	6	2
Dog	1	3

<u>Receiver Operating Characteristics (ROC), Area Under the Curve (AUC)</u> https://towardsdatascience.com/understanding-auc-roc-curve-68b2303cc9c5

- Definitions from target identification in radar receivers dyring WW2
- Operational separation choices (*threshold values*) in binary classification (*Positive* or *Negative* Prediction) depending on the administrator preference to select among *Receiver Operating Characteristics* (ROC) points
- Graph **ROC**: Function FPR \rightarrow TPR, { $0 \le$ FPR ≤ 1 , $0 \le$ TPR ≤ 1 }
- Good operational choice in **ROC**: TPR >> FPR
- Ideal choice: TPR = 1, FPR = 0

Area under ROC: AUC (Area Under the Curve) for $0 \le FPR \le 1$

Measure of *Separation Capability* of a Classifier

- No separation capability: AUC = 0.5
- Separation excellence: $AUC \gg 0.5$



Probabilistic Classification Models, Bayesian MLE & MAP Estimation (1/3)

(T. Mitchell 2016, "Machine Learning" http://www.cs.cmu.edu/~tom/mlbook/Joint_MLE_MAP.pdf)



Bayes Parametric Estimation Model

- Assumption: Elements (examples) $\mathbf{x}(i) \in \{X\}$ of the sample-space are distributed according to a *known* distribution (e.g. *Gauss*) into *discrete classes* C based on discrete probabilities of a parameter θ of the sample-space, e.g. mean of $\mathbf{x}(i)$ for every class
- Estimators $\hat{\theta}$ of parameter θ indicate the *class* C of an example $\mathbf{x}(i)$
- The $\hat{\theta}$ are inferred from observing examples $\mathbf{x}(i)$ of a **training sample** \mathcal{D} , a subset of the sample-space $\{X\}$ with outputs or **labels** $d(i) = \theta$ known to a supervisor (**teacher**) (**supervised learning**)
- According to the *Bayes Rule*:

$$P(\boldsymbol{\theta}|\boldsymbol{\mathcal{D}}) = \frac{P(\boldsymbol{\mathcal{D}}|\boldsymbol{\theta})P(\boldsymbol{\theta})}{P(\boldsymbol{\mathcal{D}})}$$

where

- $P(\mathbf{D})$: Evidence, joint probability of observing examples $\mathbf{x}(i)$ in \mathbf{D} , sampled from $\{X\}$
- $P(\theta)$: **Prior** probability of parameter θ assigned to examples in the sample-space $\{X\}$
- $P(\mathbf{D} \mid \theta)$: Likelihood, joint probability of all examples in \mathbf{D} conditioned to parameter θ
- $P(\theta \mid \mathbf{D})$: **Posterior** probability of parameter θ for all examples in \mathbf{D}

Probabilistic Classification Models, Bayesian MLE & MAP Estimation (2/3)

http://www.cs.cmu.edu/~tom/mlbook/Joint_MLE_MAP.pdf



 $P(\theta)$: *Prior Assumption* based on user-experience of $\{X\}$ *or* inferred from the *observed labeled training subset* $\mathcal{D} \subset \{X\}$ used to provide *evidence* of sample-space statistics

A Binary Parametric Classifier

Assumption: Examples x_n are *Gaussian* with average $\theta = \mu_1$ if $x_n \to C_1$ or $\theta = \mu_2$ if $x_n \to C_2$



Likelihood Densities: $f_X(x|\mathcal{C}_1), f_X(x|\mathcal{C}_2)$

Classification Errors: Shaded areas

1. Maximum Likelihood Estimation (MLE): Select θ that makes $\boldsymbol{\mathcal{D}}$ the most probable option $\hat{\theta} = \arg \max_{\boldsymbol{\theta}} P(\boldsymbol{\mathcal{D}}|\boldsymbol{\theta})$

2. Maximum a Posteriori Probability (MAP) Estimation : Select the most probable θ given the observed data \boldsymbol{D} and our *prior assumptions* summarized by $P(\theta)$

$$\hat{\theta} = \arg \max_{\theta} P(\theta | \mathbf{D}) = \arg \max_{\theta} \frac{P(\mathbf{D} | \theta) P(\theta)}{P(\mathbf{D})} = \arg \max_{\theta} P(\mathbf{D} | \theta) P(\theta)$$

Probabilistic Classification Models, Bayesian MLE & MAP Estimation (3/3)

(T. Mitchell 2016, "Machine Learning" <u>http://www.cs.cmu.edu/~tom/mlbook/Joint_MLE_MAP.pdf</u>) Example: Bernoulli experiment for tossing a coin with output $X \in \{heads, tails\} \triangleq \{1,0\}$

- Training Sample $\mathcal{D} = \{x(1), x(2), ..., x(50)\}$ of 50 *trials* (examples) used to infer $\hat{\theta}$ as an estimate of the *heads* probability $\theta = P(X = 1)$
- If the trials resulted into $a_1 = 24$ heads, $a_0 = 26$ tails, the MLE estimate of θ is $\hat{\theta} = \arg \max P(\mathcal{D}|\theta) = \frac{a_1}{(1-\theta)^{a_0}} = 0.48$, with $P(\mathcal{D}|\theta) = \theta^{a_1}(1-\theta)^{a_0}$

• For MAP estimation the *Priors*
$$P(\theta)$$
 are needed, e.g. from empirical understanding of the sampling process. In case we believe that the coin favors *heads* with $P(1) \approx 0.6$ we may modify the experiment adding to observed a_1, a_2 some *biased* imaginary examples β_1, β_2 :

$$a_1 \rightarrow a_1 + \beta_1 = 24 + 9 = 33, a_0 \rightarrow a_0 + \beta_0 = 26 + 1 = 27 \text{ and } \hat{\theta} = \frac{33}{33 + 27} = 0.55 \approx 0.6$$

• With equiprobable choices of θ , MAP \equiv MLE



Probabilistic Classification Models, Bayesian MLE & MAP Estimation (3/3)

(T. Mitchell 2016, "Machine Learning" http://www.cs.cmu.edu/~tom/mlbook/Joint MLE MAP.pdf)

Probabilities ~ Relative Frequency of examples $\{\mathbf{x}(i), d(i)\}$ in the Training Set $\boldsymbol{\mathcal{D}}$

- Input $\mathbf{x}(i) = (Gender, HoursWorked)$ of m = 2 binary coordinates (features)
- Output (label) $d(i) \cong y(i) = h(\mathbf{x}(i)) = Wealth \in \{\text{poor, rich}\}$

Gender	HoursWorked	Wealth	probability	
female	< 40.5	poor	0.2531	Estimates of Probabilities
female	< 40.5	rich	0.0246	$P(\mathbf{x}, \mathbf{v}) = P(\mathbf{G}, \mathbf{HW}, \mathbf{v})$
female	≥ 40.5	poor	0.0422	
female	≥ 40.5	rich	0.0116	where
male	< 40.5	poor	0.3313	$G \in \{M, F\}$
male	< 40.5	rich	0.0972	$HW \in \{light, hard\}$
male	≥ 40.5	poor	0.1341	$y \in \{\text{poor}, \text{rich}\}$
male	≥ 40.5	rich	0.1059	
			0.0246	

Posterior $P(y|\mathbf{x})$: $P(\operatorname{rich}|F, light) = \frac{0.0246}{0.2531+0.0246} \sim 0.09$

Gender (G)	HrsWorked (HW)	<i>P</i> (rich G,HW)	<i>P</i> (poor G,HW)
F	<40.5 (<i>light</i>)	0.09	0.91
F	>40.5 (hard)	0.21	0.79
М	<40.5 (<i>light</i>)	0.23	0.77
М	>40.5 (hard)	0.38	0.62

m = 2 features {G, HW} require 4 posterior probabilities (m features require 2^m)

Conditional Independence Assumption, Naïve Bayes Classifier (1/2)

(T. Mitchell 2016, "Machine Learning" http://www.cs.cmu.edu/~tom/mlbook/NBayesLogReg.pdf)

Bayes rule for *Random Variables X*, *Y*: $P(Y|X) = \frac{P(X|Y)}{P(X)} = \frac{P(X|Y)P(Y)}{P(X)}$

Conditional Independence of *Random Variable* $\{X|Y,Z\}$ from Y: P(X|Y,Z) = P(X|Z)

Simplifying Approximation - Naive Bayes Classifier

• For a two-dimensional sample with elements $\mathbf{x} = [x_1 \ x_2]^T$ assume that x_i 's are *Conditionally Independent* from the *output* y: $P(x_1|x_2, y) \cong P(x_1|y)$. Then:

$$P(\mathbf{x}|y) = P(x_1, x_2|y) = P(x_1|x_2, y) \times P(x_2|y) \cong P(x_1|y) \times P(x_2|y)$$

• Generalizing to *m* features of sample elements $\mathbf{x} = [x_1 \ x_2 \ ... \ x_m]^T$ the likelihood is:

$$P(\mathbf{x}|y) = P(x_1, x_2, \dots, x_m | y) \cong \prod_{k=1}^{m} P(x_k | y)$$

• The *Naive Bayes Classifier* is based on estimating the *posterior* $P(d|\mathbf{x}) \cong P(y|\mathbf{x})$ based on the *training sample likelihoods* $P(\mathbf{x}|y)$:

$$P(y|\mathbf{x}) = \frac{P(y)P(\mathbf{x}|y)}{P(\mathbf{x})} \propto P(y)P(x_1|y)P(x_2|y) \dots P(x_m|y)$$

• A new **test** sample point $\mathbf{x}^{new} = [x_1^{new} x_2^{new} \dots x_m^{new}]^T$ requires $\sim m$ likelihoods for classification instead of 2^m , a significant simplification to counter the **curse of dimensionality**

The *priors* P(y) are approximated from the occurrence frequency in the training set, assuming a *Multinomial Distribution* for **discrete** x_i **values**, or *Gaussian Distribution* (*Gaussian Naive Bayes Classifier*) for continuous x_i

Conditional Independence Assumption, Naïve Bayes Classifier (2/2)

(T. Mitchell 2016, "Machine Learning" http://www.cs.cmu.edu/~tom/mlbook/NBayesLogReg.pdf

The *Naive Bayes Classifier* is based on approximating the *posterior* $P(d|\mathbf{x}) \cong P(y|\mathbf{x})$ from the product of *conditionally independent* feature *likelihoods*

 $P(y|\mathbf{x}) \propto P(y)P(x_1|y)P(x_2|y) \dots P(x_m|y)$

Naive Bayes Algorithm:

From a *labeled training sample* $\mathbf{D} = \{(\mathbf{x}(1), d(1)), \dots, (\mathbf{x}(N), d(N))\}$ estimate:

- The priors $P(d) \cong P(y)$ for all possible classes d, e.g. $d \in \{0,1\}$ for binary classification
- The *likelihood* $P(x_k = l | y = d) \triangleq \theta_{kld}$ for every *discrete* feature x_k , k = 1, 2, ..., m of the *training elements* **x** that were classified within class (*label*) d

A new example of the *test* sample $\mathbf{x}^{new} = [x_1^{new} x_2^{new} \dots x_m^{new}]^T$, $x_k^{new} = l$ will be assigned to *class* y^{new} that will result from:

$$y^{new} \leftarrow \arg \max_{y} P(y) \prod_{k=1}^{m} P(x_k^{new} | y)$$

or
$$y^{new} \leftarrow \arg \max_{d} P(d) \prod_{k=1}^{m} \theta_{kld}$$

STOCHASTIC PROCESSES & OPTIMIZATION IN MACHINE LEARNING Example of a Naïve Bayes Classifier

https://towardsdatascience.com/all-about-naive-bayes-8e13cef044cf

Training Sample of 1000 Training Examples

Fruit	Long	Sweet	Yellow	Total
Banana	400 (80%)	350 (70%)	450 (90%)	500 (50%)
Orange	0 (0%)	150 (50%)	300 (100%)	300 (30%)
Other	100 (50%)	150 (75%)	<mark>50</mark> (25%)	200 (20%)
Total	500 (50%)	650 (65%)	800 (80%)	1000

 $P(y|\mathbf{x}) \propto P(y)P(x_1|y)P(x_2|y) \dots P(x_m|y)$

- $P(\text{Banana}|\text{Long, Sweet, Yellow}) \propto (0.5) \times (0.8) \times (0.7) \times (0.9) = 0.252$
- *P*(Orange|Long, Sweet, Yellow) = 0
- $P(\text{Other}|\text{Long, Sweet, Yellow}) \propto (0.2) \times (0.5) \times (0.75) \times (0.25) = 0.01875$

Classification of a New (Test) Example

Fruits with features (characteristics) $\mathbf{x} = (\text{Long, Sweet, Yellow})$ are classified as y = (Banana) with the greatest posterior probability $P(y|\mathbf{x}) \propto 0.252$

Note: The value of the posterior can be determined by normalizing to 1 the sum of probabilities $\frac{0.252}{0.252} = -0.021$

$$P(y|\mathbf{x}) = \frac{0.232}{0.252 + 0.01875} = 0.931$$