

A SUMMARY OF MACHINE LEARNING ALGORITHMS ¹

	K-nearest neighbour	Naive Bayes	Logistic regression
DESCRIPTION	The label of a new point \hat{x} is classified with the most frequent label \hat{t} of the k nearest training instances.	Learn $p(C_k x)$ by modelling $p(x C_k)$ and $p(C_k)$, using Bayes' rule to infer the class conditional probability. Assumes each feature \mathbf{x} independent of all others, thus 'Naive.' Suitable for data of high dimension.	Model the posterior probabilities of K classes via linear functions in x , while ensuring they fall within $[0, 1]$.
MODEL	$\hat{t} = \arg \max_C \sum_{i: x_i \in N_k(\mathbf{x}, \hat{x})} \delta(t_i, C)$ <ul style="list-style-type: none"> $N_k(\mathbf{x}, \hat{x}) \leftarrow k$ points in \mathbf{x} closest to \hat{x} Euclidean distance formula: $\sqrt{\sum_{i=1}^D (x_i - \hat{x}_i)^2}$ $\delta(a, b) \leftarrow 1$ if $a = b$; 0 o/w 	$y(\mathbf{x}) = \arg \max_k p(C_k x)$ $= \arg \max_k p(x C_k) \times p(C_k)$ $= \arg \max_k \prod_{i=1}^D p(x_i C_k) \times p(C_k)$ $= \arg \max_k \sum_{i=1}^D \log p(x_i C_k) + \log p(C_k)$	$\log \frac{P(G=1 X=x)}{P(G=K X=x)} = \beta_{10} + \beta_1^T x$ $\log \frac{P(G=2 X=x)}{P(G=K X=x)} = \beta_{20} + \beta_2^T x$ \vdots $\log \frac{P(G=K-1 X=x)}{P(G=K X=x)} = \beta_{(K-1)0} + \beta_{K-1}^T x$
OBJECTIVE	No optimization needed.	No optimisation needed.	$\hat{\uparrow}(\theta) = \sum_{i=1}^N \log p_{g_i}(x_i; \theta),$ <p>where</p> $p_k(x_i; \theta) = P(G=k X=x_i; \theta).$ <p>In two-class cases, the log-likelihood can be written as</p> $l(\beta) = \sum_{i=1}^N \{y_i \log p(x_i; \beta) + (1 - y_i) \log(1 - p(x_i; \beta))\}$ $= \sum_{i=1}^N \{y_i \beta^T x_i - \log(1 + e^{\beta^T x_i})\}.$
TRAINING	Use cross-validation to learn the appropriate k ; otherwise no training, classification based on existing points.	<p>Multivariate likelihood $p(x C_k) = \sum_{i=1}^D \log p(x_i C_k)$</p> $p_{\text{MLE}}(x_i = v C_k) = \frac{\sum_{j=1}^N \delta(t_j = C_k \wedge x_{ji} = v)}{\sum_{j=1}^N \delta(t_j = C_k)}$ <p>Multinomial likelihood $p(x C_k) = \prod_{i=1}^D p(\text{word}_i C_k)^{x_i}$</p> $p_{\text{MLE}}(\text{word}_i = v C_k) = \frac{\sum_{j=1}^N \delta(t_j = C_k) \times x_{ji}}{\sum_{j=1}^N \sum_{d=1}^D \delta(t_j = C_k) \times x_{di}},$ <p>where</p> <ul style="list-style-type: none"> x_{ji} is the count of word i in test example j; x_{di} is the count of feature d in test example j. <p>Gaussian likelihood $p(x C_k) = \prod_{i=1}^D \mathcal{N}(v; \mu_{ik}, \sigma_{ik})$</p>	<p>MLE: To maximize the log-likelihood, we set its derivatives to zero.</p> $\frac{\partial l(\beta)}{\partial \beta} = \sum_{i=1}^N x_i (y_i - p(x_i; \beta)) = 0.$ <p>Gradient descent is written as:</p> $\theta_{k+1} = \theta_k - \eta_k g_k,$ <p>where η is the learning rate. The selection of η and other training methods including Newton's method should refer to [?].</p>

¹Initiated by Emanuel Ferm, and enriched by Guanqun Cao (guanqun.cao@tut.fi) on July 23, 2013.

	K-nearest neighbour	Naive Bayes	Logistic regression
REGULARIZATION	k acts as to regularise the classifier: as $k \rightarrow N$ the boundary becomes smoother.	Use a Dirichlet prior on the parameters to obtain a MAP estimate.	The L_1 penalty used in the lasso for variable selection and shrinkage applies for any linear regression model. Here we have
		Multivariate likelihood $p_{\text{MAP}}(x_i = v \mathcal{C}_k) = \frac{(\beta_i - 1) + \sum_{j=1}^N \delta(t_j = \mathcal{C}_k \wedge x_{ji} = v)}{ x_i (\beta_i - 1) + \sum_{j=1}^N \delta(t_j = \mathcal{C}_k)}$	$\arg \max \left\{ \sum_{i=1}^N [y_i(\beta_0 + \beta^T x_i) - \log(1 + e^{\beta_0 + \beta^T x_i} - \lambda \sum_{j=1}^p \ \beta_j\)] \right\}.$
		Multinomial likelihood $p_{\text{MAP}}(\text{word}_i = v \mathcal{C}_k) = \frac{(\alpha_i - 1) + \sum_{j=1}^N \delta(t_j = \mathcal{C}_k) \times x_{ji}}{\sum_{j=1}^N \sum_{d=1}^D (\delta(t_j = \mathcal{C}_k) \times x_{di}) - D + \sum_{d=1}^D \alpha_d}$	
COMPLEXITY	$\mathcal{O}(NM)$ space complexity, since all training instances and all their features need to be kept in memory.	$\mathcal{O}(NM)$, each training instance must be visited and each of its features counted.	$\mathcal{O}(INMK)$, since each training instance must be visited and each combination of class and features must be calculated for the appropriate feature mapping.
NON-LINEAR	Natively finds non-linear boundaries.	Can only learn linear boundaries for multivariate/multinomial attributes. With Gaussian attributes, quadratic boundaries can be learned with uni-modal distributions.	Kernel smoothing methods can estimate the regression function $f(x)$ over the domain R^p by fitting a different but simple model separately at each query point x_0 , (not to be confused with kernel methods). This localization is achieved via a weighting function or kernel $K_\lambda(x_0, x_i)$, which assigns a weight to x_i based on its distance from x_0 . The kernels K_λ are typically indexed by a parameter λ that dictates the width of the neighborhood. Locally weighted regression solves a separate weighted least squares problem at each target point x_0 :
			$\min_{\alpha(x_0), \beta(x_0)} \sum_{i=1}^N K_\lambda(x_0, x_i) [y_i - \alpha(x_0) - \beta(x_0)x_i]^2.$
ONLINE LEARNING	To be added.	To be added.	The objective in online learning is the <i>regret</i> , which is the averaged loss incurred relative to the best using a single fixed parameter value:
			$\text{regret}_k = \frac{1}{k} \sum_{t=1}^k f(\theta_t, z_t) - \min_{\theta^* \in \Theta} \frac{1}{k} \sum_{t=1}^k f(\theta^*, z_t),$
			where $f(\theta, z_t)$ is the loss function of certain kind. <i>Online gradient descent</i> describes that at each step k , update the parameters using
			$\theta_{k+1} = \text{proj}_\Theta(\theta_k - \eta_k g_k),$
			where $\text{proj}_V(v) = \arg \min_{w \in V} \ w - v\ _2$ is the projection of vector v onto space \mathcal{V} , $g_k = \nabla f(\theta_k, z_k)$ is the gradient, and η_k is the step size.

	Perceptron	Support vector machines	Linear discriminant analysis
DESCRIPTION	Directly estimate the linear function $y(x)$ by iteratively updating the weight vector when incorrectly classifying a training instance.	A maximum margin classifier: finds the separating hyperplane with the maximum margin to its closest data points.	
MODEL	<p>Binary, linear classifier:</p> $y(x) = \text{sign}(\mathbf{w}^T x),$ <p>where:</p> $\text{sign}(x) = \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}$ <p>Multiclass perceptron:</p> $y(x) = \arg \max_{C_k} \mathbf{w}^T \phi(x, C_k)$	$y(x) = \sum_{n=1}^N \lambda_n t_n x^T x_n + w_0$	
OBJECTIVE	<p>Tries to minimise the Error function; the number of incorrectly classified input vectors:</p> $\arg \min_{\mathbf{w}} E_P(\mathbf{w}) = \arg \min_{\mathbf{w}} - \sum_{n \in \mathcal{M}} \mathbf{w}^T x_n t_n,$ <p>where \mathcal{M} is the set of misclassified training vectors.</p>	<p>Primal</p> $\arg \min_{\mathbf{w}, w_0} \frac{1}{2} \ \mathbf{w}\ ^2$ <p>s.t. $t_n(\mathbf{w}^T x_n + w_0) \geq 1 \quad \forall n$</p> <p>Dual</p> $\tilde{\mathcal{L}}(\wedge) = \sum_{n=1}^N \lambda_n - \sum_{n=1}^N \sum_{m=1}^N \lambda_n \lambda_m t_n t_m x_n^T x_m$ <p>s.t. $\lambda_n \geq 0, \quad \sum_{n=1}^N \lambda_n t_n = 0, \quad \forall n$</p>	
TRAINING	<p>Iterate over each training example x_n, and update the weight vector if misclassification:</p> $\begin{aligned} \mathbf{w}^{i+1} &= \mathbf{w}^i + \eta \Delta E_P(\mathbf{w}) \\ &= \mathbf{w}^i + \eta x_n t_n, \end{aligned}$ <p>where typically $\eta = 1$.</p> <p>For the multiclass perceptron:</p> $\mathbf{w}^{i+1} = \mathbf{w}^i + \phi(x, t) - \phi(x, y(x))$	<ul style="list-style-type: none"> • Quadratic Programming (QP) • SMO, Sequential Minimal Optimisation (chunking). 	

Perceptron

REGULARIZATION

The Voted Perceptron: run the perceptron i times and store each iteration's weight vector. Then:

$$y(x) = \text{sign} \left(\sum_i c_i \times \text{sign}(\mathbf{w}_i^T x) \right),$$

where c_i is the number of correctly classified training instances for \mathbf{w}_i .

COMPLEXITY

$\mathcal{O}(INML)$, since each combination of instance, class and features must be calculated.

NON-LINEAR

Use a kernel $K(x, x')$, and 1 weight per training instance:

$$y(x) = \text{sign} \left(\sum_{n=1}^N w_n t_n K(x, x_n) \right)$$

... and the update:

$$w_n^{i+1} = w_n^i + 1$$

ONLINE LEARNING

The perceptron is an online algorithm per default.

Support vector machines

The soft margin SVM: penalise a hyperplane by the number and distance of misclassified points.

Primal

$$\arg \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n$$

$$\text{s.t. } t_n(\mathbf{w}^T x_n + w_0) \geq 1 - \xi_n, \quad \xi_n > 0 \quad \forall n$$

Dual

$$\tilde{\mathcal{L}}(\lambda) = \sum_{n=1}^N \lambda_n - \sum_{n=1}^N \sum_{m=1}^N \lambda_n \lambda_m t_n t_m x_n^T x_m$$

$$\text{s.t. } 0 \leq \lambda_n \leq C, \quad \sum_{n=1}^N \lambda_n t_n = 0, \quad \forall n$$

- QP: $\mathcal{O}(n^3)$;
- SMO: much more efficient than QP, since computation based only on support vectors.

Use a non-linear kernel $K(x, x')$:

$$\begin{aligned} y(x) &= \sum_{n=1}^N \lambda_n t_n x^T x_n + w_0 \\ &= \sum_{n=1}^N \lambda_n t_n K(x, x_n) + w_0 \end{aligned}$$

$$\begin{aligned} \tilde{\mathcal{L}}(\lambda) &= \sum_{n=1}^N \lambda_n - \sum_{n=1}^N \sum_{m=1}^N \lambda_n \lambda_m t_n t_m x_n^T x_m \\ &= \sum_{n=1}^N \lambda_n - \sum_{n=1}^N \sum_{m=1}^N \lambda_n \lambda_m t_n t_m K(x_n, x_m) \end{aligned}$$

Online SVM. See, for example:

- *The Huller: A Simple and Efficient Online SVM*, Bordes & Bottou (2005)
- *Pegasos: Primal Estimated sub-Gradient Solver for SVM*, Shalev-Shwartz et al. (2007)

Linear discriminant analysis

Mixture Models

DESCRIPTION A probabilistic clustering algorithm, where clusters are modelled as latent Gaussians and each data point is assigned the probability of being drawn from a particular Gaussian.

MODEL Assignments to clusters by specifying probabilities

$$p(x^{(i)}, z^{(i)}) = p(x^{(i)}|z^{(i)})p(z^{(i)})$$

...with $z^{(i)} \sim \text{Multinomial}(\gamma)$, and $\gamma_{nk} \equiv p(k|x_n)$ s.t. $\sum_{j=1}^K \gamma_{nj} = 1$. I.e. want to maximise the probability of the observed data \mathbf{x} .

OBJECTIVE

$$\begin{aligned} \mathcal{L}(\mathbf{x}, \pi, \mu, \Sigma) &= \log p(\mathbf{x}|\pi, \mu, \Sigma) \\ &= \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right) \end{aligned}$$

TRAINING **Expectation:** For each n, k set:

$$\begin{aligned} \gamma_{nk} &= p(z^{(i)} = k | x^{(i)}; \gamma, \mu, \Sigma) \quad (= p(k|x_n)) \\ &= \frac{p(x^{(i)}|z^{(i)} = k; \mu, \Sigma)p(z^{(i)} = k; \pi)}{\sum_{j=1}^K p(x^{(i)}|z^{(i)} = j; \mu, \Sigma)p(z^{(i)} = j; \pi)} \\ &= \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)} \end{aligned}$$

Maximisation:

$$\begin{aligned} \pi_k &= \frac{1}{N} \sum_{n=1}^N \gamma_{nk} \\ \Sigma_k &= \frac{\sum_{n=1}^N \gamma_{nk} (x_n - \mu_k)(x_n - \mu_k)^T}{\sum_{n=1}^N \gamma_{nk}} \\ \mu_k &= \frac{\sum_{n=1}^N \gamma_{nk} x_n}{\sum_{n=1}^N \gamma_{nk}} \end{aligned}$$

k-means

A hard-margin, geometric clustering algorithm, where each data point is assigned to its closest centroid.

Hard assignments $r_{nk} \in \{0, 1\}$ s.t. $\forall n \sum_k r_{nk} = 1$, i.e. each data point is assigned to one cluster k .

Geometric distance: The Euclidean distance, l^2 norm:

$$\|x_n - \mu_k\|_2 = \sqrt{\sum_{i=1}^D (x_{ni} - \mu_{ki})^2}$$

$$\arg \min_{\mathbf{r}, \mu} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|_2^2$$

... i.e. minimise the distance from each cluster centre to each of its points.

Expectation:

$$r_{nk} = \begin{cases} 1 & \text{if } \|x_n - \mu_k\|^2 \text{ minimal for } k \\ 0 & \text{o/w} \end{cases}$$

Maximisation:

$$\mu_{\text{MLE}}^{(k)} = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}},$$

where $\mu^{(k)}$ is the centroid of cluster k .

	Mixture Models	<i>k</i> -means
REGULARIZATION	The mixture of Gaussians assigns probabilities for each cluster to each data point, and as such is capable of capturing ambiguities in the data set.	Only hard-margin assignment to clusters.
COMPLEXITY	To be added.	The heuristic solution under the assumption of each point perturbed by a normal distribution with mean 0 and variance σ^2 incurs a time complexity $O(n^{34}k^{34}d^8\log^4(n)/\sigma^6)$.
NON-LINEAR	Not applicable.	For non-linearly separable data, use kernel <i>k</i> -means as suggested in: <i>Kernel k-means, Spectral Clustering and Normalized Cuts</i> , Dhillon et al. (2004).
ONLINE LEARNING	Online Gaussian Mixture Models. A good start is: <i>A View of the EM Algorithm that Justifies Incremental, Sparse, and Other Variants</i> , Neal & Hinton (1998).	Sequential <i>k</i> -means: update the centroids after processing one point at a time.

References

- [1] Kevin P Murphy. *Machine learning: a probabilistic perspective*. The MIT Press, 2012.
- [2] Trevor. Hastie, Robert. Tibshirani, and J Jerome H Friedman. *The elements of statistical learning*, volume 1. Springer New York, (Online version), 2013.
- [3] Richard O Duda, Peter E Hart, and David G Stork. *Pattern classification*. John Wiley & Sons, 2012.