

Equations

Bayes' rule	$P(A B) = P(A)P(B A)/P(B)$
Expected value	$\mathbb{E}[cX] = c\mathbb{E}[X]$
Variance	$Var(X) = E[(X - E[X])^2] = \sigma^2$ $Var(X) = E[X^2] - (E[X])^2$ $Var(cX) = c^2 Var(X)$ $\mathbb{E}[(x_i - \mu_i)(x_j - \mu_j)]$ $X^T X = U S^2 U^T$ for centered dataset
Covariance	$\sigma_{ij} = \frac{1}{N} \sum_{k=1:N} (x_i^k - \mu_i)(x_j^k - \mu_j)$
Covariance matrix	$[\partial^2 R / \partial w_i \partial w_j]$
Hessian	$\partial^2 R / \partial w^2$
Hessian matrix	$\ x - y\ _p = (\sum_{j=1}^d x_j - y_j ^p)^{1/p}$
p-norm	

Distributions

Isotropic normal	$\frac{1}{\sigma\sqrt{2\pi}} \exp(-(x - \mu)^2 / (2\sigma^2))$
Multivariate normal	$\frac{1}{\sqrt{(2\pi)^k \Sigma }} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu))$
Poisson distribution	$p(x, \lambda) = (e^{-\lambda} \lambda^x) / x!$

Matrices

$(AB)^T$	$B^T A^T$
$(A^T)^{-1}$	$(A^{-1})^T$
$\partial(x^T a) / \partial x$	a
$\partial(x^T A x) / \partial x$	$(A + A^T)x$
$\partial \text{Trace}(XA) / \partial X$	A^T
$\partial(a^T X b) / \partial X$	ab^T

Miscellaneous Notes

Posterior $P(Y = i|x)$

Techniques

Maximum Likelihood Estimation Find likelihood function, take negative log likelihood, take derivative with respect to optimized parameter. Set derivative to 0. $\text{argmax}[P(\text{data}|\text{model})]$

Maximum A Posteriori (MAP)
 $\text{argmin}[-\log P(D|f)(-\log P(f))]$. $-\log P(f)$ serves the role of regularizer. $\text{argmax}[P(\text{data}|\text{model})P(\text{model})]$

Lagrange Multiplier To maximize $f(x, y)$ subject to constraint $g(x, y) = 0$, examine the equation using a Lagrange multiplier: $\mathcal{L}(x, y, \lambda) = f(x, y) + \lambda g(x, y)$. Then, take the derivative of the equation with respect to x, y , and λ .

Linear Models

Centroid method Find class centroids, subtract one centroid from the other to get weight vector, construct separating hyperplane with $f(x) = w \cdot x = 0$. Equivalent to Hebb's rule with equal number of points.

Linear SVM Fits a hyperplane to separate and classify data. Soft-margin hyperparameter C represents tradeoff between robustness and fit. Minimizes $1/\text{margin} + C(\text{training error})$. Approaches hard margin as $C \rightarrow \infty$, and vice versa. Classify new points via $\text{sign}(\sum_{i=1}^n \alpha_i K(x^{(i)}, x))$.

Hebb's rule $w = \sum_k y^k x^k$, $w \leftarrow w + y^k \Phi(x^k)$, $\alpha^k \leftarrow \alpha^k + y^k$

Perceptron algorithm Always finds separating hyperplane if samples linearly separable. Equations: $\Delta w_i = \eta y \Phi_i(x_i)$, $\Delta \alpha_k = \eta y^k$. Normal perceptron: $z < 0$, update if misclassified. Large-margin: $z < 1$, update if misclassified or within margin. Optimum margin: $\min(z)$, only update if closest point to margin.

Newton's Method $w \leftarrow w - H^{-1} \nabla_w R$

Stochastic Gradient Descent $\Delta w_i = -\eta \partial L / \partial w_i$. Obtain dual form $\Delta \alpha$ by using kernel trick, NOT by differentiating loss. Does not exploit convexity of risk function; best approach for big data but NOT when N or d is small.

Update rules:

least mean squares: $\Delta w_i = \eta(y - f(x))\Phi_i(x)$
 $\Delta \alpha_k = \eta(y^k - f(x^k))$

Kernels

Measure similarity between two datapoints. Represents a dot product in some (potentially infinite) feature space: $k(s, t) = \Phi(s) \cdot \Phi(t)$ (don't need to know Φ representation). Good kernels are symmetric: $k(x, x') = k(x', x)$. Kernel matrix should be invertible, possibly after regularization $(K + \lambda I)$. Satisfied if matrix is PSD, all eigenvalues $\lambda_i > 0$.

Kernel Trick

parametric	$f(x) = w \cdot \Phi(x)$ $w = \sum_k \alpha_k \Phi(x^k)$
non-parametric	$f(x) = \sum_k \alpha_k k(x^k, x)$ $k(x^k, x) = \Phi(x^k) \cdot \Phi(x)$ $f(x) = \sum_{k=1:N} y_k k(x, x_k)$

Parzen window

Radial kernels

Top-hat $k(x, x_k) = 1(\|x - x_k\|^2 < \sigma^2)$

Gaussian $\exp(-\|x - x_k\|^2 / 2\sigma^2)$

Non-radial kernels

Linear $k(x, x_k) = x \cdot x_k$

Polynomial $k(x, x_k) = (1 + x \cdot x_k)^q$

Risk Minimization

As model complexity increases, bias decreases while variance increases.

Loss functions

Square	$(1 - z)^2$
Perceptron	$\max(0, -z)$
Hinge	$\max(0, 1 - z)$
Logistic	$\log(1 + e^{-z})$
Adaboost	e^{-z}

Risk/Loss Functionals

functional margin	$z = yf(x), y = \pm 1$
risk	$(1/N) \sum_{k=1:N} L(f(x^k), y^k)$
batch gradient	$w \leftarrow w - \eta \nabla_w R$
stochastic gradient	$w \leftarrow w - \eta \nabla_w L$
SRM/regularization	$w_i \leftarrow (1 - \gamma)w_i - \eta \partial R_{train} / \partial w_i$
linear discriminant	$f(x) = \sum_i w_i x_i = wx$

Regularization penalizes model complexity at expense of more training error, often explicitly part of loss function. Takes the form of shrinkage (weight decay, ridge regression, SVM), feature selection, kernel parameters, etc.

Risk Types

guaranteed risk: $R_{\text{gua}}[f] = R_{\text{train}}[f] + \epsilon(\delta, C/N)$ with high probability $(1 - \delta)$

regularized risk: $R_{\text{reg}}[f] = R_{\text{train}}[f] + \lambda \|w\|^2$

Norms $\|x\|_p = (|x_1|^p + |x_2|^p + \dots)^{1/p}$

L_0 norm	penalizes number of features considered
L_1 norm (lasso)	makes weight vector more sparse
L_2 norm (ridge)	shrinks weight vector to reduce variance

Logistic Regression

Like Hebb's rule but weighted; misclassifications are more heavily weighted (multiplied by $S(-z)$). A link function links the functional margin z and the likelihood $P(D|f)$.

Linear logistic regression log odds ratio (logit) is a linear function of x , $\ln[P_f(Y = 1|X = x)/P_f(Y = -1|X = x)] = w \cdot x + b$

Logistic function $S(t) = g^{-1}(t) = 1/(1 + e^{-t})$
Training finds the optimal β value such that if $P(Y = 1|x) > 0.5$, output 1; otherwise, output 0.

$$P(Y = 1|x) = 1/(1 + \exp(-\beta^T x)) \quad (1)$$

$$\begin{aligned} \Delta w_i &= -\eta \partial L / \partial w_i = -\eta \partial L / \partial z \cdot \partial z / \partial w_i \\ &= \eta S(-z) y x_i \end{aligned}$$

Update equations with shrinkage:

$$\Delta w_i = (-\gamma w_i) + \eta S(-z) y \Phi_i(x)$$

$$\Delta \alpha_k = (-\gamma \alpha_k) + \eta S(-z) y^k \text{ for example } k$$

$$\Delta \alpha_h = (-\gamma \alpha_h) \text{ for other examples}$$

Ridge Regression

$$\sum_i w_i x_i^k = y^k \text{ for all } k = 1..m$$

$$Xw^T = y$$

$$X^T X w^T = X^T y$$

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

Optimal solution: $w^T = X^+ y$

Residual $y - \hat{y} = (I - X X^+) y$

Pseudo-inverse

Case 1) $N > d$ overdetermined, no exact solution. Optimal RSS solution is

$$X^+ = \lim_{\lambda \rightarrow 0} (X^T X + \lambda I)^{-1} X^T$$

Case 2) $N < d$ underdetermined, optimize for $\min(\|w\|)$

$$X^+ = \lim_{\lambda \rightarrow 0} X^T (X X^T + \lambda I)^{-1}$$

Not limit when $\lambda \rightarrow 0$, but find optimal value through cross-validation.

Kernel trick In case 2 where $N \ll d$, dimensionality of features can approach ∞ . Instead, replace XX^T by a (N, N) kernel matrix $K = k(x^k, x^h)$. $\alpha = (K + \lambda I)^{-1} y$ yields the nonlinear regression function $f(x) = \sum_k \alpha_k k(x, x_k)$.

Principal Component Analysis (PCA) decrease dimensionality of features by constructing linear combinations of the features such that the reconstructed patterns are as close as possible to the original features (minimize RSS). Accomplishes this by removing the dimensions with the smallest eigenvalues (smallest variance, affects the data the least).

PCA Solution: $X' = XU$, $X'' = X'U^T = XU^T$, $\min_U \|X - XU U^T\|^2$

Bayesian Decision Theory

Optimal decision boundary $P(Y = 1|x) = P(Y = -1|x)$

Assumptions Datasets should be IID for naive Bayes.

Bonferroni Correction: $p' = mp$, where we use m classifiers.

Performance Evaluation

Metrics error rate, bit-error rate (BER), cross-entropy, area under the ROC curve (AUC), R^2

Error bar Σ of the error rate is $\sqrt{E(1 - E)/n}$

Sample size Number of test examples needed is inversely proportional to error rate. Rule of thumb: $n = 100/E$

Model Selection

Filters reduce number of hyperparameters by measuring feature relevance. Does not make use of the learning machine, uses statistics instead. Techniques used are prior knowledge, unsupervised learning to reduce dimensionality of input space, transfer learning, meta-learning, surrogate learning machines

Wrappers use the learning machine to evaluate the performance of alternative feature subsets. Use a search method to explore the space of possible subsets. Considered a black box, no knowledge of the learning algorithm is needed to apply this method.

Embedded methods reduce the number of hyperparameters as much as possible before passing into wrapper method (too many hyperparameters leads to overfitting). Pushes hyperparameters to a lower level and have them estimated in the process of learning parameters. This cannot be done by optimizing the training error; need two levels of inference in order to enjoy finite capacity of the learning problem. Optimizing the kernel parameters (kernel with in the Gaussian kernel or degree in the polynomial kernel) leads to infinite VC dimension (can learn perfectly any training set); optimizing the ridge or regularization parameter Λ leads to zero capacity ($w = 0$ is the optimum).

Search strategies grid search, simulated annealing, pattern search

Voting methods multiple models cast a vote for a class, classifier weighs each vote and makes a decision.

Gaussian Classification

Assumptions: Generating model (draw y first, draw x given y), variance in dataset explained by Gaussian noise, independence of features in given class (no covariance), same variance for all classes. NOT optimum Bayes classifier because assumptions almost always violated. We can post-fit bias term by adjusting the threshold. If two classes have same variance, Gaussian classifier is a linear discriminant (equivalent to centroid method, Hebb's rule with target values $y^k = 1/N_1$ if in class 1 and $-1/N_0$ otherwise). Posterior of Gaussian model is logistic with Poisson or Gaussian class conditionals.

Bayes' rule $P(Y = y|X = x) = P(Y = y)P(X = x|Y = y)$

Prior $P(Y = y)$, relative class abundance (occurrences over total count)

Likelihood $P(X = x|Y = y)$, probability x belongs to class y , proportional to $\exp(-\|x - \mu^{[y]}\|^2/2\sigma^2)$, Gaussian noise model

Linear Discriminant Analysis

Centering: subtracting mean of the features

Standardizing sphering; subtracting by mean, dividing by standard deviation (component-wise), $\Phi(x^k) = (x^k - \mu)/\sigma$

Principal axes Diagonalize $\Sigma = X^T X = U S^2 U^T$, where U is an orthogonal matrix of column eigenvectors such that $U^T U = I$ and S^2 is a diagonal matrix of positive eigenvalues s_i^2 . Matrix

$R = U S U^T = \Sigma^{1/2}$. To obtain a diagonal covariance matrix S , $\Phi = X U$. To axis-align the Gaussian, apply $U^{-1}(X - \mu)$

Whitening multiply by square root of inverse covariance matrix, $\Phi = X(X^T X)^{-1/2} = X R^{-1} = X \Sigma^{-1/2}$

Linear Discriminant Analysis Generalization of Gaussian classifier for cases where the input variables aren't statistically independent, but all classes have same covariance matrix Σ . Once the input space is rotated into the principal axes of the covariance matrix and rescaled by the eigenvalues, LDA is like the isotropic Gaussian classifier (centroid method). Useful for multi-class classification and data visualization.

Rotate input space, then find the direction that maximizes the ratio of the between class variance over the within class variance.

$$\max \left(\frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2} \right) \quad \text{test} \quad (3)$$

PCA, ridge regression use covariance matrix of all data combined. LCA uses pooled, within-class covariance. In both cases, whitening is applied with $\Phi = X \Sigma^{-1/2}$.

Dimensionality Reduction

Singular Value Decomposition

Decomposition of $X = U S V^T$. $X^T X = V S S V^T = V D V^T$, $X X^T = U S S U^T = U D U^T$. Left singular vectors of X are eigenvectors of $X X^T$. Right singular vectors are eigenvectors of $X^T X$. Non-zero singular values are the square roots of the nonzero eigenvalues of both $X X^T$ and $X^T X$. Calculate eigenvectors of $X^T X$ and put eigenvectors in rows for right singular matrix. To get left singular vectors (arranged in columns), use $u_i = X v_i / \sigma_i$ where σ_i is the square root of the eigenvalue.

Principal Component Analysis

Subtract mean from each point, (optionally) scale each dimension by its variance. Compute covariance matrix $C = X^T X$, then find the k largest eigenvalues of $C = V D V^T$. Eigenvalue magnitude indicates variance along that axis.

Decision Trees

Bagging Generate bootstrap replicates of original training data by sampling with replacement.

Random Forests Ensemble method by creating multiple trees via bagging. Additional restriction that at each node, only a subset of splits is allowed. Each tree votes for a class, and the mode of votes

is taken as the result. Used to lower variance on low-bias, high-variance random decision trees.

Boosting Weight weak learners according to how well they perform on certain parts of the dataset. Points that were incorrectly classified are weighted more heavily on future update iterations. Creates a strong learner from a group of weak learners. Primarily reduces bias. Weights initialized to 1, then updated as $w_t \leftarrow w_t \exp(-y_t H_t)$.

AdaBoost The bigger the error ϵ , the smaller the hypothesis weight multiplier α (the weights adapt to the error). Boosts examples if incorrectly predicted.

Neural Networks

Back-Propagation

$$\frac{\partial J}{\partial w_{ij}^{(l)}} = \frac{\partial J(w)}{\partial s_j^{(l)}} \cdot \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}} = \frac{\partial J(w)}{\partial s_j^{(l)}} \cdot x^{(l-1)} \quad (2)$$

Nonlinearities

Sigmoid $s(t) = 1/(1 + e^{-t}) \in [0, 1]$, $ds(t)/dt = s(t)(1 - s(t))$. Gradients can become small for values away from 0, leading to vanishing gradient problem.

tanh $\tanh(t) \in [-1, 1]$

Convolutional Neural Nets

Used to combat curse of dimensionality, doesn't need as many weights as a normal neural net.

Convolution Connectivity is local in space but full in depth. Number of parameters for each convolutional layer equals filter size times number of filters.

Pooling Makes representations smaller without losing too much information. Does not have any weights (uses a function like max-pooling, L2 pooling, etc.).

Nonlinearities functions are applied on every pixel of each layer after the pooling stage.

Clustering

k-nn error is bounded by twice the Bayes risk.

k-means Global minimization is NP-hard.

k-medoids

Gaussian Mixture Models

Nonparametric Density Estimation

Histogramming: Construct histogram from feature space by creating bins and counting points.

Kernel Density Estimation: Histogram might have lots of empty bins; create Gaussians at every point, and combine them to create a density estimation with Parzen windows.