

# Contents

<b>Contents</b> .....	i		
<b>Notation</b> .....	v		
<b>1 Introduction</b> .....	1		
1.1 Types of machine learning .....	1		
1.2 Supervised learning .....	1		
1.2.1 Classification .....	1		
1.2.2 Regression .....	2		
1.3 Unsupervised learning .....	2		
1.3.1 Discovering clusters .....	2		
1.3.2 Discovering latent factors .....	2		
1.3.3 Discovering graph structure .....	2		
1.3.4 Matrix completion .....	2		
1.4 Three elements of a machine learning model .....	2		
1.4.1 Representation .....	2		
1.4.2 Evaluation .....	3		
1.4.3 Optimization .....	4		
1.5 Model Selection .....	4		
1.6 The Curse of Dimensionality .....	5		
1.7 Decision Theory .....	5		
1.7.1 Minimizing the misclassification rate .....	6		
1.7.2 Minimizing the expected loss .....	6		
1.7.3 The reject option .....	6		
1.7.4 Inference and decision .....	6		
1.8 Information theory .....	7		
1.8.1 Relative entropy and mutual information .....	9		
1.9 Some basic concepts .....	10		
1.9.1 Parametric vs non-parametric models .....	10		
1.9.2 A simple non-parametric classifier: K-nearest neighbours .....	10		
1.9.3 Overfitting .....	11		
1.9.4 Model selection .....	11		
<b>2 Probability and Statistics</b> .....	13		
2.1 Frequentists vs. Bayesians .....	13		
2.2 probability theory .....	13		
2.2.1 Concepts .....	13		
2.2.2 Fundamental rules .....	14		
2.2.3 Multivariate random variables .....	15		
2.2.4 Bayes rule .....	16		
2.2.5 Independence and conditional independence .....	17		
2.2.6 Quantiles .....	17		
		2.2.7 Mean and variance .....	18
		2.3 Some common discrete distributions .....	18
		2.3.1 The Bernoulli and binomial distributions .....	19
		2.3.2 The multinoulli and multinomial distributions .....	19
		2.3.3 The Poisson distribution .....	19
		2.3.4 The empirical distribution .....	20
		2.4 Some common continuous distributions .....	20
		2.4.1 Gaussian (normal) distribution .....	20
		2.4.2 Student's t-distribution .....	21
		2.4.3 The Laplace distribution .....	21
		2.4.4 The gamma distribution .....	22
		2.4.5 The beta distribution .....	23
		2.4.6 Pareto distribution .....	25
		2.5 The Gaussian Distribution .....	26
		2.5.1 Conditional Gaussian distributions .....	26
		2.5.2 Marginal Gaussian distributions .....	26
		2.5.3 Baye's theorem for Gaussian variables .....	26
		2.6 The exponential family .....	26
		2.6.1 Definition .....	26
		2.6.2 Maximum likelihood and sufficient statistics .....	27
		2.6.3 Conjugate priors .....	27
		2.6.4 Noninformative priors .....	27
		2.6.5 Examples .....	28
		2.6.6 Log partition function .....	29
		2.6.7 MLE for the exponential family .....	30
		2.6.8 Bayes for the exponential family .....	31
		2.6.9 Maximum entropy derivation of the exponential family * .....	31
		2.7 Nonparametric Methods .....	31
		2.7.1 Kernel density estimators .....	31
		2.7.2 Nearest-neighbour methods .....	31
		2.8 Joint probability distributions .....	31
		2.8.1 Covariance and correlation .....	31
		2.8.2 Multivariate Gaussian distribution .....	32
		2.8.3 Multivariate Student's t-distribution .....	33
		2.8.4 Dirichlet distribution .....	33
		2.9 Transformations of random variables .....	36
		2.9.1 Linear transformations .....	36
		2.9.2 General transformations .....	37

2.9.3	Central limit theorem	37	4.3	Inference in jointly Gaussian distributions	59
2.10	Monte Carlo approximation	38	4.3.1	Statement of the result	59
2.11	Information theory	38	4.3.2	Examples	60
2.11.1	Entropy	38	4.3.3	information form	60
2.11.2	KL divergence	38	4.3.4	Proof	60
2.11.3	Mutual information	39	4.4	Linear Gaussian systems	60
<b>3</b>	<b>Generative models for discrete data</b>	41	4.4.1	Statement of the result	61
3.1	Generative classifier	41	4.4.2	Examples	61
3.2	Bayesian concept learning	41	4.4.3	Proof	61
3.2.1	Likelihood	41	4.5	Digression: The Wishart distribution *	61
3.2.2	Prior	41	4.6	Inferring the parameters of an MVN	61
3.2.3	Posterior	41	4.6.1	Posterior distribution of $\mu$	61
3.2.4	Posterior predictive distribution	42	4.6.2	Posterior distribution of $\Sigma$ *	61
3.3	The beta-binomial model	42	4.6.3	Posterior distribution of $\mu$ and $\Sigma$ *	61
3.3.1	Likelihood	42	4.6.4	Sensor fusion with unknown precisions *	61
3.3.2	Prior	42	<b>5</b>	<b>Bayesian statistics</b>	63
3.3.3	Posterior	43	5.1	Introduction	63
3.3.4	Posterior predictive distribution	44	5.2	Summarizing posterior distributions	63
3.4	The Dirichlet-multinomial model	45	5.2.1	MAP estimation	63
3.4.1	Likelihood	45	5.2.2	Credible intervals	65
3.4.2	Prior	45	5.2.3	Inference for a difference in proportions	65
3.4.3	Posterior	46	5.3	Bayesian model selection	66
3.4.4	Posterior predictive distribution	46	5.3.1	Bayesian Occam's razor	66
3.5	Naive Bayes classifiers	46	5.3.2	Computing the marginal likelihood (evidence)	67
3.5.1	Optimization	47	5.3.3	Bayes factors	70
3.5.2	Using the model for prediction	48	5.4	Priors	70
3.5.3	The log-sum-exp trick	48	5.4.1	Uninformative priors	70
3.5.4	Feature selection using mutual information	48	5.4.2	Robust priors	70
3.5.5	Classifying documents using bag of words	49	5.4.3	Mixtures of conjugate priors	70
<b>4</b>	<b>Gaussian Models</b>	51	5.5	Hierarchical Bayes	71
4.1	Basics	51	5.6	Empirical Bayes	71
4.1.1	proof	52	5.7	Bayesian decision theory	71
4.1.2	limitations	53	5.7.1	Bayes estimators for common loss functions	72
4.1.3	MLE for a MVN	54	5.7.2	The false positive vs false negative tradeoff	73
4.1.4	Maximum entropy derivation of the Gaussian *	54	<b>6</b>	<b>Linear Regression</b>	75
4.2	Gaussian discriminant analysis	54	6.1	Introduction	75
4.2.1	Quadratic discriminant analysis (QDA)	54	6.2	Representation	75
4.2.2	Linear discriminant analysis (LDA)	55	6.3	Maximum likelihood estimations(least squares)	76
4.2.3	Two-class LDA	57	6.3.1	Derivations of the MLE	77
4.2.4	MLE for discriminant analysis	58	6.3.2	Geometric interpretation	79
4.2.5	Strategies for preventing overfitting	59	6.3.3	Sequential learning	80
4.2.6	Regularized LDA *	59	6.4	Ridge regression(MAP)	81
4.2.7	Diagonal LDA	59	6.4.1	Basic idea	81
4.2.8	Nearest shrunken centroids classifier *	59	6.4.2	Multiple outputs	82

6.4.3	Numerically stable computation *	82	8.2.4	Gradient descent optimization	108
6.4.4	Connection with PCA *	82	8.3	Error Backpropagation	108
6.4.5	Regularization effects of big data	82	8.3.1	Evaluation of error-function derivatives	108
6.5	The Bias-Variance Decomposition	83	8.3.2	Efficiency of back propagation	110
6.6	Bayesian linear regression	85	8.3.3	The Jacobian matrix	110
6.6.1	Parameter distribution	85	8.4	The Hessian Matrix	110
6.6.2	Predictive distribution	85	8.5	Regularization in Neural Networks	110
6.6.3	Equivalent kernel	86	8.6	Mixture Density Networks	110
6.7	Bayesian Model Comparison	87	8.7	Bayesian Neural Networks	110
<b>7</b>	<b>Linear Models for Classification</b>	89	<b>9</b>	<b>Sparse Kernel Machines</b>	111
7.1	Discriminant Functions	89	9.1	Introduction	111
7.1.1	Two classes	89	9.2	Maximum Margin Classifiers	111
7.1.2	Multiple classes	90	9.2.1	Overlapping class distributions	113
7.1.3	Least squares for classification	90	9.3	Relevance Vector Machines	118
7.1.4	Fisher's linear discriminant	90	<b>10</b>	<b>Mixture models and EM</b>	119
7.1.5	Extensions to higher dimensions and multiple classes	92	10.1	Introduction	119
7.1.6	Probabilistic interpretation of FLDA *	93	10.2	K-means Clustering	119
7.1.7	Relation to least squares	93	10.2.1	representation	119
7.1.8	Fisher's discriminant for multiple classes	93	10.2.2	evaluation	119
7.1.9	The perceptron algorithm	93	10.2.3	optimization	119
7.2	Probabilistic Generative Models	94	10.3	Mixtures of Gaussians	120
7.2.1	Two classes	94	10.3.1	representation	120
7.2.2	Multiple classes	95	10.3.2	Maximum likelihood	121
7.2.3	Continuous inputs	95	10.3.3	EM for Gaussian mixtures	122
7.2.4	Maximum likelihood solution	96	10.4	Latent Variables View of EM	124
7.2.5	Discrete features	97	10.4.1	representation	124
7.2.6	Exponential family	98	10.5	The EM Algorithm in General	124
7.3	Probabilistic Discriminative Models	98	10.5.1	Maximum posterior	126
7.3.1	Fixed basis functions	98	10.5.2	generalized EM	127
7.3.2	Logistic regression	99	<b>11</b>	<b>Continuous Latent Variables</b>	129
7.3.3	Iterative reweighted least squares	100	11.1	Principal Component Analysis	129
7.3.4	Multiclass logistic regression	101	11.1.1	Introduction	129
7.3.5	Probit regression	102	11.1.2	Maximum variance formulation	129
7.3.6	Canonical link functions	102	11.1.3	Minimum-error formulation	130
7.3.7	The Laplace Approximation	103	11.1.4	Applications of PCA	130
7.3.8	Model comparison and BIC	103	11.1.5	PCA for high-dimensional data	130
7.4	Bayesian Logistic Regression	103	11.2	Probabilistic PCA	130
7.4.1	Laplace approximation	103	11.3	Kernel PCA	130
7.4.2	Predictive distribution	103	11.4	Nonlinear Latent Variable Models	130
<b>8</b>	<b>Neural Networks</b>	105	<b>12</b>	<b>Combining Models</b>	131
8.1	Feed-forward Network Functions	105	12.1	Introduction	131
8.1.1	Weight-space symmetries	106	12.2	Bayesian Model Averaging	131
8.2	Network Training	106	12.3	Committees	131
8.2.1	Parameter optimization	107	12.4	Boosting	131
8.2.2	Local quadratic approximation	108	12.4.1	Minimizing exponential error	132
8.2.3	Use of gradient information	108	12.4.2	Error functions for boosting	134
			12.5	Tree-based Models	134
			12.6	Conditional Mixture Models	134

<b>A</b>	<b>Matrix Decomposition</b>	135				
A.1	LU Decomposition	135				
A.2	QR Decomposition	135				
A.3	Cholesky Decomposition	135				
A.4	Eigen value Decomposition	135				
A.5	Singular Value Decomposition	135				
	A.5.1 Definition	135				
	A.5.2 Proof	135				
	A.5.3 Application	136				
<b>B</b>	<b>Optimization methods</b>	137				
B.1	Convexity	137				
B.2	Gradient descent	137				
	B.2.1 Stochastic gradient descent	137				
	B.2.2 Batch gradient descent	138				
	B.2.3 Line search	138				
	B.2.4 Momentum term	138				
B.3	Lagrange duality	138				
	B.3.1 Primal form	138				
	B.3.2 Dual form	138				
B.4	Newton's method	138				
B.5	Quasi-Newton method	139				
	B.5.1 DFP	139				
	B.5.2 BFGS	139				
	B.5.3 Broyden	140				
	<b>Glossary</b>	141				

# Notation

## Introduction

It is very difficult to come up with a single, consistent notation to cover the wide variety of data, models and algorithms that we discuss. Furthermore, conventions differ between machine learning and statistics, and between different books and papers. Nevertheless, we have tried to be as consistent as possible. Below we summarize most of the notation used in this book, although individual sections may introduce new notation. Note also that the same symbol may have different meanings depending on the context, although we try to avoid this where possible.

## General math notation

Symbol	Meaning
$\lfloor x \rfloor$	Floor of $x$ , i.e., round down to nearest integer
$\lceil x \rceil$	Ceiling of $x$ , i.e., round up to nearest integer
$\mathbf{x} \otimes \mathbf{y}$	Convolution of $\mathbf{x}$ and $\mathbf{y}$
$\mathbf{x} \odot \mathbf{y}$	Hadamard (elementwise) product of $\mathbf{x}$ and $\mathbf{y}$
$a \wedge b$	logical AND
$a \vee b$	logical OR
$\neg a$	logical NOT
$\mathbb{I}(x)$	Indicator function, $\mathbb{I}(x) = 1$ if $x$ is true, else $\mathbb{I}(x) = 0$
$\infty$	Infinity
$\rightarrow$	Tends towards, e.g., $n \rightarrow \infty$
$\propto$	Proportional to, so $y = ax$ can be written as $y \propto x$
$ x $	Absolute value
$ \mathcal{S} $	Size (cardinality) of a set
$n!$	Factorial function
$\nabla$	Vector of first derivatives
$\nabla^2$	Hessian matrix of second derivatives
$\triangleq$	Defined as
$O(\cdot)$	Big-O: roughly means order of magnitude
$\mathbb{R}$	The real numbers
$1:n$	Range (Matlab convention): $1:n = 1, 2, \dots, n$
$\approx$	Approximately equal to
$\arg \max_x f(x)$	Argmax: the value $x$ that maximizes $f$
$B(a, b)$	Beta function, $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$
$B(\alpha)$	Multivariate beta function, $\frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum_k \alpha_k)}$
$\binom{n}{k}$	$n$ choose $k$ , equal to $n!/(k!(n-k)!)$
$\delta(x)$	Dirac delta function, $\delta(x) = \infty$ if $x = 0$ , else $\delta(x) = 0$
$\exp(x)$	Exponential function $e^x$
$\Gamma(x)$	Gamma function, $\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$
$\Psi(x)$	Digamma function, $\Psi(x) = \frac{d}{dx} \log \Gamma(x)$

---

$\mathcal{X}$	A set from which values are drawn (e.g., $\mathcal{X} = \mathbb{R}^D$ )
---------------	---

---

## Linear algebra notation

We use boldface lower-case to denote vectors, such as  $\mathbf{x}$ , and boldface upper-case to denote matrices, such as  $\mathbf{X}$ . We denote entries in a matrix by non-bold upper case letters, such as  $X_{ij}$ .

Vectors are assumed to be column vectors, unless noted otherwise. We use  $(x_1, \dots, x_D)$  to denote a column vector created by stacking  $D$  scalars. If we write  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ , where the left hand side is a matrix, we mean to stack the  $\mathbf{x}_i$  along the columns, creating a matrix.

Symbol	Meaning
$\mathbf{X} \succ 0$	$\mathbf{X}$ is a positive definite matrix
$tr(\mathbf{X})$	Trace of a matrix
$det(\mathbf{X})$	Determinant of matrix $\mathbf{X}$
$ \mathbf{X} $	Determinant of matrix $\mathbf{X}$
$\mathbf{X}^{-1}$	Inverse of a matrix
$\mathbf{X}^\dagger$	Pseudo-inverse of a matrix
$\mathbf{X}^T$	Transpose of a matrix
$\mathbf{x}^T$	Transpose of a vector
$diag(\mathbf{x})$	Diagonal matrix made from vector $\mathbf{x}$
$diag(\mathbf{X})$	Diagonal vector extracted from matrix $\mathbf{X}$
$\mathbf{I}$ or $\mathbf{I}_d$	Identity matrix of size $d \times d$ (ones on diagonal, zeros of)
$\mathbf{1}$ or $\mathbf{1}_d$	Vector of ones (of length $d$ )
$\mathbf{0}$ or $\mathbf{0}_d$	Vector of zeros (of length $d$ )
$\ \mathbf{x}\  = \ \mathbf{x}\ _2$	Euclidean or $\ell_2$ norm $\sqrt{\sum_{j=1}^d x_j^2}$
$\ \mathbf{x}\ _1$	$\ell_1$ norm $\sum_{j=1}^d  x_j $
$\mathbf{X}_{:,j}$	$j$ 'th column of matrix
$\mathbf{X}_{i,:}$	transpose of $i$ 'th row of matrix (a column vector)
$\mathbf{X}_{i,j}$	Element $(i, j)$ of matrix $\mathbf{X}$
$\mathbf{x} \otimes \mathbf{y}$	Tensor product of $\mathbf{x}$ and $\mathbf{y}$

---

## Probability notation

We denote random and fixed scalars by lower case, random and fixed vectors by bold lower case, and random and fixed matrices by bold upper case. Occasionally we use non-bold upper case to denote scalar random variables. Also, we use  $p()$  for both discrete and continuous random variables

Symbol	Meaning
$X, Y$	Random variable
$P()$	Probability of a random event
$F()$	Cumulative distribution function(CDF), also called distribution function
$p(x)$	Probability mass function(PMF)
$f(x)$	probability density function(PDF)
$F(x, y)$	Joint CDF
$p(x, y)$	Joint PMF
$f(x, y)$	Joint PDF

$p(X Y)$	Conditional PMF, also called conditional probability
$f_{X Y}(x y)$	Conditional PDF
$X \perp Y$	X is independent of Y
$X \not\perp Y$	X is not independent of Y
$X \perp Y Z$	X is conditionally independent of Y given Z
$X \not\perp Y Z$	X is not conditionally independent of Y given Z
$X \sim p$	X is distributed according to distribution $p$
$\alpha$	Parameters of a Beta or Dirichlet distribution
$\text{cov}[X]$	Covariance of X
$\mathbb{E}[X]$	Expected value of X
$\mathbb{E}_q[X]$	Expected value of X wrt distribution $q$
$\mathbb{H}(X)$ or $\mathbb{H}(p)$	Entropy of distribution $p(X)$
$\mathbb{I}(X;Y)$	Mutual information between X and Y
$\mathbb{KL}(p  q)$	KL divergence from distribution $p$ to $q$
$\ell(\theta)$	Log-likelihood function
$L(\theta, a)$	Loss function for taking action $a$ when true state of nature is $\theta$
$\lambda$	Precision (inverse variance) $\lambda = 1/\sigma^2$
$\Lambda$	Precision matrix $\Lambda = \Sigma^{-1}$
$\text{mode}[\mathbf{X}]$	Most probable value of $\mathbf{X}$
$\mu$	Mean of a scalar distribution
$\boldsymbol{\mu}$	Mean of a multivariate distribution
$\Phi$	cdf of standard normal
$\phi$	pdf of standard normal
$\pi$	multinomial parameter vector, Stationary distribution of Markov chain
$\rho$	Correlation coefficient
$\text{sigm}(x)$	Sigmoid (logistic) function, $\frac{1}{1 + e^{-x}}$
$\sigma^2$	Variance
$\Sigma$	Covariance matrix
$\text{var}[x]$	Variance of $x$
$\nu$	Degrees of freedom parameter
$Z$	Normalization constant of a probability distribution

## Machine learning/statistics notation

In general, we use upper case letters to denote constants, such as  $C, K, M, N, T$ , etc. We use lower case letters as dummy indexes of the appropriate range, such as  $c = 1 : C$  to index classes,  $i = 1 : M$  to index data cases,  $j = 1 : N$  to index input features,  $k = 1 : K$  to index states or clusters,  $t = 1 : T$  to index time, etc.

We use  $x$  to represent an observed data vector. In a supervised problem, we use  $y$  or  $\mathbf{y}$  to represent the desired output label. We use  $z$  to represent a hidden variable. Sometimes we also use  $q$  to represent a hidden discrete variable.

Symbol	Meaning
$C$	Number of classes
$D$	Dimensionality of data vector (number of features)
$N$	Number of data cases
$N_c$	Number of examples of class $c$ , $N_c = \sum_{i=1}^N \mathbb{I}(y_i = c)$
$R$	Number of outputs (response variables)
$\mathcal{D}$	Training data $\mathcal{D} = \{(\mathbf{x}_i, y_i)   i = 1 : N\}$
$\mathcal{D}_{test}$	Test data
$\mathcal{X}$	Input space
$\mathcal{Y}$	Output space

$K$	Number of states or dimensions of a variable (often latent)
$k(x, y)$	Kernel function
$\mathbf{K}$	Kernel matrix
$\mathcal{H}$	Hypothesis space
$L$	Loss function
$J(\boldsymbol{\theta})$	Cost function
$f(\mathbf{x})$	Decision function
$P(y \mathbf{x})$	TODO
$\lambda$	Strength of $\ell_2$ or $\ell_1$ <i>regularizer</i>
$\phi(x)$	Basis function expansion of feature vector $\mathbf{x}$
$\Phi$	Basis function expansion of design matrix $\mathbf{X}$
$q()$	Approximate or proposal distribution
$Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{old})$	Auxiliary function in EM
$T$	Length of a sequence
$T(\mathcal{D})$	Test statistic for data
$\mathbf{T}$	Transition matrix of Markov chain
$\boldsymbol{\theta}$	Parameter vector
$\boldsymbol{\theta}^{(s)}$	$s$ 'th sample of parameter vector
$\hat{\boldsymbol{\theta}}$	Estimate (usually MLE or MAP) of $\boldsymbol{\theta}$
$\hat{\boldsymbol{\theta}}_{MLE}$	Maximum likelihood estimate of $\boldsymbol{\theta}$
$\hat{\boldsymbol{\theta}}_{MAP}$	MAP estimate of $\boldsymbol{\theta}$
$\bar{\boldsymbol{\theta}}$	Estimate (usually posterior mean) of $\boldsymbol{\theta}$
$\mathbf{w}$	Vector of regression weights (called $\boldsymbol{\beta}$ in statistics)
$b$	intercept (called $\varepsilon$ in statistics)
$\mathbf{W}$	Matrix of regression weights
$x_{ij}$	Component (i.e., feature) $j$ of data case $i$ , for $i = 1 : N, j = 1 : D$
$\mathbf{x}_i$	Training case, $i = 1 : N$
$\mathbf{X}$	Design matrix of size $N \times D$
$\bar{\mathbf{x}}$	Empirical mean $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$
$\tilde{\mathbf{x}}$	Future test case
$\mathbf{x}_*$	Feature test case
$\mathbf{y}$	Vector of all training labels $\mathbf{y} = (y_1, \dots, y_N)$
$z_{ij}$	Latent component $j$ for case $i$

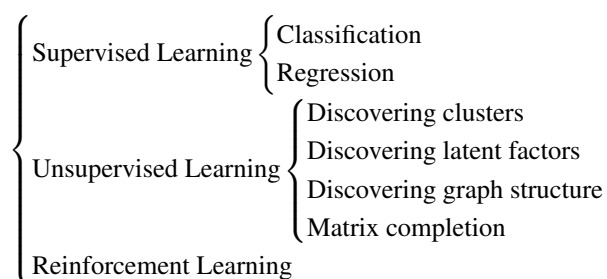
---



# Chapter 1

## Introduction

### 1.1 Types of machine learning



In the **predictive** or **supervised learning** approach, the goal is to learn a **mapping** from **inputs  $\mathbf{x}$**  to **outputs  $y$** , given a labeled set of input-output pairs  $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ . Here  $D$  is called the **training set**, and  $N$  is the number of training examples. In the simplest setting, each training input  $\mathbf{x}_i$  is a  $D$ -dimensional vector of numbers, representing, say, the height and weight of a person, which are called **features, attributes, or covariates**.

### 1.2 Supervised learning

Similarly the form of the output or **response variable** can in principle be anything, but most methods assume that  $y_i$  is **categorical or nominal** variable from some finite set,  $y_i \in \{1, \dots, C\}$ . When  $y_i$  is categorical, the problem is known as **classification or pattern recognition**, and when real-valued, known as **regression**.

#### 1.2.1 Classification

Here the goal is to learn a mapping from inputs  $x$  to outputs  $y$ , where  $y \in \{1, \dots, C\}$ , with  $C$  being the number of classes. If  $C = 2$ , this is called **binary classification**; if  $C > 2$ , this is called **multiclass classification**. If the class labels are not mutually exclusive, we call it **multi-label classification**, but this is best viewed as predicting multiple related binary class labels (a so-called **multiple output model**). One way to formalize the problem is as **function approximation**: assume  $y = f(\mathbf{x})$  for some unknown function  $f$ , and the goal of learning is to estimate the function  $f$  given a labeled training set, and then to make predictions (estimate) using  $\hat{y} = \hat{f}(\mathbf{x})$ . Our main goal is to make predictions on novel inputs, meaning ones that we have not seen before (**generalization**).

##### 1.2.1.1 Probabilistic predictions

We denote the probability distribution over possible labels, given the input vector  $\mathbf{x}$  and training set  $\mathcal{D}$  by  $p(y|\mathbf{x}, \mathcal{D})$ . Given a probabilistic output, we can always compute out "best guess" as to the "true label" using

$$\hat{y} = \hat{f}(\mathbf{x}) = \arg \max_{c=1}^C p(y = c|\mathbf{x}, \mathcal{D}) \quad (1.1)$$

This corresponds to the most probable class label, and is called the **mode** of distribution  $p(y|\mathbf{x}, \mathcal{D})$ ; it is also known as a **MAP estimate** (MAP stands for **maximum a posteriori**).

### 1.2.1.2 Applications

## 1.2.2 Regression

## 1.3 Unsupervised learning

**Descriptive or unsupervised learning** approach is sometimes called **knowledge discovery**. We will formalize our task as one of **density estimation**, that is we want to build models of the form  $p(\mathbf{x}_i|\theta)$ , instead of  $p(y_i|\mathbf{x}_i, \theta)$ .

### 1.3.1 Discovering clusters

Let  $z_i \in \{1, \dots, K\}$  represent the cluster to which data point  $i$  is assigned. ( $z_i$  is an example of **hidden or latent** variable).

### 1.3.2 Discovering latent factors

Although the data may appear high dimensional, there may only be a small number of degrees of variability, corresponding to **latent factors**. The most common approach to dimensionality reduction is called **principal components analysis** or **PCA**.

### 1.3.3 Discovering graph structure

### 1.3.4 Matrix completion

#### 1.3.4.1 Image inpainting

#### 1.3.4.2 Collaborative filtering

#### 1.3.4.3 Market basket analysis

## 1.4 Three elements of a machine learning model

**Model = Representation + Evaluation + Optimization**<sup>1</sup>

### 1.4.1 Representation

In supervised learning, a model must be represented as a conditional probability distribution  $P(y|x)$  (usually we call it classifier) or a decision function  $f(x)$ . The set of classifiers (or decision functions) is called the hypothesis space of the model. Choosing a representation for a model is tantamount to choosing the hypothesis space that it can possibly learn.

---

<sup>1</sup> Domingos, P. A few useful things to know about machine learning. Commun. ACM. 55(10):7887 (2012).

## 1.4.2 Evaluation

In the hypothesis space, an evaluation function (also called objective function or risk function) is needed to distinguish good classifiers (or decision functions) from bad ones.

### 1.4.2.1 Loss function and risk function

**Definition 1.1.** In order to measure how well a function fits the training data, a **loss function**  $L : Y \times Y \rightarrow R \geq 0$  is defined. For training example  $(x_i, y_i)$ , the loss of predicting the value  $\hat{y}$  is  $L(y_i, \hat{y})$ .

The following is some common loss functions:

1. 0-1 loss function

$$L(Y, f(X)) = \mathbb{I}(Y, f(X)) = \begin{cases} 1, & Y \neq f(X) \\ 0, & Y = f(X) \end{cases}$$

2. Quadratic(squared) loss function  $L(Y, f(X)) = \frac{1}{2} (Y - f(X))^2$

3. Absolute loss function  $L(Y, f(X)) = |Y - f(X)|$

4. Exponential loss function  $L(Y, f(X)) = \exp(-\hat{y}_i f(\mathbf{x}_i))$

5. Logarithmic loss function

$$L(Y, P(Y|X)) = -\log P(Y|X)$$

Name	Loss	Derivative	$f^*$	Algorithm
Squared error	$\frac{1}{2} (y_i - f(\mathbf{x}_i))^2$	$y_i - f(\mathbf{x}_i)$	$\mathbb{E}[y \mathbf{x}_i]$	L2Boosting
Absolute error	$ y_i - f(\mathbf{x}_i) $	$\text{sgn}(y_i - f(\mathbf{x}_i))$	$\text{median}(y \mathbf{x}_i)$	Gradient boosting
Exponential loss	$\exp(-\hat{y}_i f(\mathbf{x}_i))$	$-\hat{y}_i \exp(-\hat{y}_i f(\mathbf{x}_i))$	$\frac{1}{2} \log \frac{\pi_i}{1-\pi_i}$	AdaBoost
Logloss	$\log(1 + e^{-\hat{y}_i f_i})$	$y_i - \pi_i$	$\frac{1}{2} \log \frac{\pi_i}{1-\pi_i}$	LogitBoost

**Definition 1.2.** The risk of function  $f$  is defined as the expected loss of  $f$ :

$$R_{\text{exp}}(f) = E[L(Y, f(X))] = \int L(y, f(x)) P(x, y) dx dy \quad (1.2)$$

which is also called expected loss or **risk function**.

**Definition 1.3.** The risk function  $R_{\text{exp}}(f)$  can be estimated from the training data as

$$R_{\text{emp}}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) \quad (1.3)$$

which is also called empirical loss or **empirical risk**.

You can define your own loss function, but if you're a novice, you're probably better off using one from the literature. There are conditions that loss functions should meet<sup>2</sup>:

1. They should approximate the actual loss you're trying to minimize. As was said in the other answer, the standard loss functions for classification is zero-one-loss (misclassification rate) and the ones used for training classifiers are approximations of that loss.
2. The loss function should work with your intended optimization algorithm. That's why zero-one-loss is not used directly: it doesn't work with gradient-based optimization methods since it doesn't have a well-defined gradient (or even a subgradient, like the hinge loss for SVMs has).

The main algorithm that optimizes the zero-one-loss directly is the old perceptron algorithm(chapter §??).

<sup>2</sup> <http://t.cn/zTrDxLO>

### 1.4.2.2 ERM and SRM

**Definition 1.4.** ERM(Empirical risk minimization)

$$\min_{f \in \mathcal{F}} R_{\text{emp}}(f) = \min_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) \quad (1.4)$$

**Definition 1.5.** Structural risk

$$R_{\text{smp}}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f) \quad (1.5)$$

**Definition 1.6.** SRM(Structural risk minimization)

$$\min_{f \in \mathcal{F}} R_{\text{srm}}(f) = \min_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f) \quad (1.6)$$

## 1.4.3 Optimization

Finally, we need a **training algorithm**(also called **learning algorithm**) to search among the classifiers in the hypothesis space for the highest-scoring one. The choice of optimization technique is key to the **efficiency** of the model.

## 1.5 Model Selection

If data is plentiful, then one approach to avoid over-fitting is to use some of the available data to train a range of models, or a given model with a range of values for its complexity parameters, and then to compare them on independent data, sometimes called a **Validation set**, and select the one having the best predictive performance. If the model design is iterated many times using a limited size data set, then some over-fitting to the validation data can occur and so it may be necessary to keep aside a third **test set** on which the performance of the selected model is finally evaluated.

### 1.5.0.1 Cross validation

**Definition 1.7. Cross validation**, sometimes called *rotation estimation*, is a *model validation* technique for assessing how the results of a statistical analysis will generalize to an independent data set<sup>3</sup>.

Common types of cross-validation:

1. K-fold cross-validation. In k-fold cross-validation, the original sample is randomly partitioned into k equal size subsamples. Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k - 1 subsamples are used as training data.
2. 2-fold cross-validation. Also, called simple cross-validation or holdout method. This is the simplest variation of k-fold cross-validation, k=2.
3. Leave-one-out cross-validation(*LOOCV*). k=M, the number of original samples.

In many applications, however, the supply of data for training and testing will be limited. We wish to use as much of available data as possible for both training and validation set. One solution to this dilemma is to use **cross-validations**. One round of cross-validation involves partitioning a sample of data into complementary subsets, performing the analysis on one subset (called the training set), and validating the analysis on the other subset (called the

<sup>3</sup> [http://en.wikipedia.org/wiki/Cross-validation\\_\(statistics\)](http://en.wikipedia.org/wiki/Cross-validation_(statistics))

validation set or testing set). To reduce variability, multiple rounds of cross-validation are performed using different partitions, and the validation results are averaged over the rounds.

## 1.6 The Curse of Dimensionality

The common theme that arise when analyzing and organizing data in **high-dimensional spaces** is that when the **dimensionality increases**, the **volume of the space** increases so fast (**grows exponentially with the dimensionality**) that the available data become sparse.

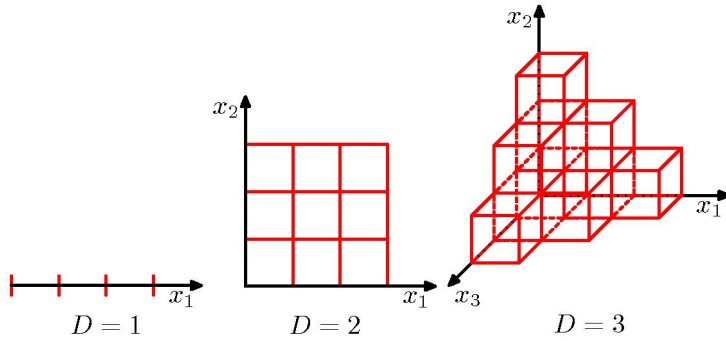


Fig. 1.1: Illustration of the curse of dimensionality: the number of regions of a regular grid grows exponentially with the dimensionality  $D$  of the space.

The volume of a sphere of radius  $r$  in  $D$  dimensions must scale as  $r^D$ , and so we write

$$V_D(r) = K_D r^D \quad (1.7)$$

The fraction of the volume of the sphere that lies between radius  $r = 1 - \epsilon$  and  $r = 1$

$$\frac{V_D(1) - V_D(1 - \epsilon)}{V_D(1)} = 1 - (1 - \epsilon)^D \quad (1.8)$$

for large  $D$  this fraction tends to 1 even for small values of  $\epsilon$ . Thus most of the volume of a sphere is concentrated in a thin shell near the surface!

## 1.7 Decision Theory

Probability theory provides us with a consistent mathematical framework for quantifying and manipulating uncertainty. And when combined with probability theory, decision theory allows us to make optimal decisions in situations involving uncertainty such as those encountered in pattern recognition.

Suppose we have an input vector  $x$  together with a corresponding vector  $t$  of target variables, and our goal is to predict  $t$  given a new value for  $x$ . For regression problems,  $t$  will comprise continuous variables, whereas for classification problems  $t$  will represent class labels. The joint probability distribution  $p(x, t)$  provides a complete summary of the uncertainty associated with these variables. Determination of  $p(x, t)$  from a set of training data is an example of **inference** and is typically a very difficult problem.

The **decision** step, is the subject of decision theory to tell us how to make optimal decisions given the appropriate probabilities. We shall see that the decision stage is generally very simple, even trivial, once we have solved the inference problem.

Consider, for a classification problem, in which case, the input are vectors  $\mathbf{x}$  and we denote the class  $i$  by  $\mathcal{C}_k$ . Using Bayes' theorem, these probabilities can be expressed in the form

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} \quad (1.9)$$

Note that any of the quantities appearing in Bayes' theorem can be obtained from the joint distribution  $p(\mathbf{x}, \mathcal{C}_k)$  by either **marginalizing** or **conditioning** with respect to the appropriate variables. We can now interpret  $p(\mathcal{C}_k)$  as the prior probability for the class  $\mathcal{C}_k$ , and the  $p(\mathcal{C}_k|\mathbf{x})$  as the corresponding posterior probability,  $p(\mathbf{x}|\mathcal{C}_k)$  is the likelihood function. If our aim is to minimize the chance of assigning  $\mathbf{x}$  to the wrong class, then intuitively we would choose the class having the higher posterior probability.

### 1.7.1 Minimizing the misclassification rate

### 1.7.2 Minimizing the expected loss

### 1.7.3 The reject option

### 1.7.4 Inference and decision

We have broken the classification problem down into two separate stages, the **inference stage** in which we use training data to learn to **model** for  $p(\mathcal{C}_k|\mathbf{x})$ , and the subsequent **decision stage** in which we use these posterior probabilities to make optimal class assignments. An alternative possibility would be to solve both problems together and simply learn a function that maps input  $\mathbf{x}$  directly into decisions. Such a function is called a **discriminant function**.

In fact, we can identify three distinct approaches to solving decision problems, given, in decreasing order of complexity, by

1. First solve the inference problem of determining the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$  for each class  $\mathcal{C}_k$  individually. Also separately infer the prior class probabilities  $p(\mathcal{C}_k)$ . Then use *Bayes* theorem in the form

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} \quad (1.10)$$

to find the posterior class probabilities  $p(\mathcal{C}_k|\mathbf{x})$ . As usual, the denominator in *Bayes* theorem can be found in terms of the quantities appearing in the numerator, because

$$p(\mathbf{x}) = \sum_k p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k) \quad (1.11)$$

Equivalently, we can model the joint distribution  $p(\mathbf{x}, \mathcal{C}_k)$  directly and then normalize to obtain the posterior probabilities. Having found the posterior probabilities, we use decision theory to determine class membership for each new input  $\mathbf{x}$ . Approaches that explicitly or implicitly model the **distribution of inputs as well as outputs** are known as **Generative models**, because by sampling from them it is possible to generate synthetic data points in the input space.

2. First solve the inference problem of determining the posterior class probabilities  $p(\mathcal{C}_k|\mathbf{x})$ , and then subsequently use decision theory to assign each new  $\mathbf{x}$  to one of the classes. Approaches that model the posterior probabilities are called **discriminative models**
3. Find a function  $f(\mathbf{x})$ , called a discriminant function, which maps each input  $\mathbf{x}$  directly into a class label. For instance, in the case of two-class problems,  $f(\cdot)$  might be binary valued and such that  $f = 0$  represents class  $\mathcal{C}_1$  and so on.

Now consider the relative merits of these three alternatives. Approach (a) is the most demanding because it involves finding the joint distribution over both  $\mathbf{x}$  and  $\mathcal{C}_k$ . For many applications,  $\mathbf{x}$  will have high dimensionality, and conse-

quently we may need a large training set in order to be able to determine the class-conditional densities to reasonable accuracy. Note that the class **priors**  $p(C_k)$  can often be estimated simply from the **fractions of the training set data points in each of the classes**. One advantage of approach (a), however, is that it also allows the marginal density of data  $p(x)$  to be determined from (1.83). This can be useful for detecting new data points that have low probability under the model and for which the predictions may be of low accuracy, which is known as **outlier detection or novelty detection** (Bishop, 1994; Tarassenko, 1995).

However, if we only wish to make classification decisions, then it can be wasteful of computational resources, and excessively demanding of data, to find the joint distribution  $p(x, C_k)$  when in fact we only really need the posterior probabilities  $p(C_k|x)$ , which can be obtained directly through approach (b). Indeed, the class conditional densities may contain a lot of structure that has little effect on the posterior probabilities, as illustrated in Figure 1.27. There has been much interest in exploring the relative merits of generative and discriminative approaches to machine learning, and in finding ways to combine them (Jebara, 2004; Lasserre et al., 2006). An even simpler approach is (c) in which we use the training data to find a discriminant function  $f(x)$  that maps each  $x$  directly onto a class label, thereby combining the inference and decision stages into a single learning problem. In the example of Figure 1.27, this would correspond to finding the value of  $x$  shown by the vertical green line, because this is the decision boundary giving the minimum probability of misclassification.

Reasons for wanting to compute the posterior probabilities include:

- **Minimize risk.**
- **Reject option.**
- **Compensating for class priors.**
- **Combining models.** For complex applications, we may wish to break the problem into a number of smaller sub-problems each of which can be tackled by a separate module. As long as each of the two models gives posterior probabilities for the classes, we can combine the outputs systematically using the rules of probability. One simple way to do this is to assume that, for each class separately, the distributions of inputs are independent, so that

$$p(x_I, x_B|C_k) = p(x_I|C_k)p(x_B|C_k) \quad (1.12)$$

This is an example of **conditional independence** property. The posterior probability, given both data (from different modules), is then given by

$$p(C_k|x_I, x_B) \propto p(x_I, x_B|C_k)p(C_k) \quad (1.13)$$

$$\propto p(x_I|C_k)p(x_B|C_k)p(C_k) \quad (1.14)$$

$$\propto \frac{p(C_k|x_I)p(C_k|x_B)}{p(C_k)} \quad (1.15)$$

## 1.8 Information theory

Consider a discrete random variable  $x$  and we ask how much information is received when we observe a specific value for this variable. Our measure of information content will therefore depend on the probability distribution  $p(x)$ , and we therefore look for a quantity  $h(x)$  that is a monotonic function of the probability  $p(x)$  and that expresses the information content. The form of  $h(\cdot)$  can be found by noting that if we have two events  $x$  and  $y$  that are unrelated, then the information gain from observing both of them should be the sum of the information gained from each of them separately, so that  $h(x, y) = h(x) + h(y)$ . Two unrelated events will be statistically independent and so  $p(x, y) = p(x)p(y)$ . From these two relationships, it is easily shown that  $h(x)$  must be given by the logarithm of  $p(x)$  and so we have

$$h(x) = -\log_2 p(x) \quad (1.16)$$

where the negative sign ensures that information is positive or zero. Note that low probability events  $x$  correspond to high information content. The choice of basis for the logarithm is arbitrary.

Now suppose that a sender wishes to transmit the value of a random variable to a receiver. The average amount of information that they transmit in the process is obtained by taking the expectation of 1.16 with respect to the distribution  $p(x)$  and is given by

$$H[x] = - \sum_x p(x) \log_2 p(x) \quad (1.17)$$

This important quantity is called the **entropy** of the random variable  $x$ .

The **noiseless coding theorem** (Shannon, 1948) states that the entropy is a lower bound on the number of bits needed to transmit the state of a random variable.

For example, we can represent the states  $x = a, b, c, d, e, f, g, h$  using, for instance, the following set of code strings: 0, 10, 110, 1110, 111100, 111101, 111110, 111111. The average length of the code that has to be transmitted is then

$$\text{average code length} = \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64} \times 6 = 2 \text{ bits} \quad (1.18)$$

Now, we switch to the use of natural logarithm in defining entropy. And entropy can be understood in an alternative view of **multiplicity**. The continuous form of the entropy is given by

$$\lim_{\Delta \rightarrow 0} \left\{ \sum_i p(x_i) \Delta \ln p(x_i) \right\} = - \int p(x) \ln p(x) dx \quad (1.19)$$

where the quantity on the right-side is called the **differential entropy**.

For a density defined over multiple continuous variables, denoted collectively by the vector  $\mathbf{x}$ , the differential entropy is given by

$$H[\mathbf{x}] = - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \quad (1.20)$$

In the case of discrete distributions, we see that the maximum entropy configuration corresponds to an equal distribution of probabilities across the possible states of variable. Now consider the maximum entropy configuration for a continuous variable. It will be necessary to constrain the first and second moments of  $p(x)$  as well as preserving the normalization constraint. Therefore we need to maximize the differential entropy  $H[x] = - \int p(x) \ln p(x) dx$  with three constraints

$$\int_{-\infty}^{\infty} p(x) dx = 1 \quad (1.21)$$

$$\int_{-\infty}^{\infty} x p(x) dx = \mu \quad (1.22)$$

$$\int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx = \sigma^2. \quad (1.23)$$

The constrained maximization can be performed using Lagrange multipliers so that we maximize the following functional with respect to  $p(x)$

$$- \int_{-\infty}^{\infty} p(x) \ln p(x) dx + \lambda_1 \left( \int_{-\infty}^{\infty} p(x) dx - 1 \right) + \lambda_2 \left( \int_{-\infty}^{\infty} x p(x) dx - \mu \right) + \lambda_3 \left( \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx - \sigma^2 \right) \quad (1.24)$$

Using the calculus of variations, we set the derivative of this functional to zero giving

$$p(x) = \exp\{-1 + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2\} \quad (1.25)$$

The Lagrange multipliers can be found by back substitution of this result into the three constraint equations, leading finally to the result

$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\} \quad (1.26)$$

and so the distribution that maximizes the differential entropy is the Gaussian.

If we evaluate the differential entropy of the Gaussian, we obtain

$$H[x] = \frac{1}{2} \{1 + \ln(2\pi\sigma^2)\} \quad (1.27)$$

Thus we see again that the entropy increases as the distribution becomes broader, i.e., as  $\sigma^2$  increases.



Suppose we have joint distribution  $p(\mathbf{x}, \mathbf{y})$  from which we draw pairs of values of  $\mathbf{x}$  and  $\mathbf{y}$ . If a value of  $\mathbf{x}$  is already known, then the additional information needed to specify the corresponding value of  $\mathbf{y}$  is given by  $-\ln p(\mathbf{y}|\mathbf{x})$ , thus the average additional information needed can be written as

$$H[\mathbf{y}|\mathbf{x}] = - \iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{y} d\mathbf{x} \quad (1.28)$$

which is called the **conditional entropy** of  $\mathbf{y}$  given  $\mathbf{x}$ . Using the product rule, the conditional entropy satisfies the relation

$$H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}] \quad (1.29)$$

where  $H[\mathbf{x}, \mathbf{y}]$  is the differential entropy of  $p(\mathbf{x}, \mathbf{y})$  and  $H[\mathbf{x}]$  is the differential entropy of the marginal distribution  $p(\mathbf{x})$ .

*Proof.*

$$H[\mathbf{x}, \mathbf{y}] - H[\mathbf{y}|\mathbf{x}] = - \iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{y}, \mathbf{x}) d\mathbf{y} d\mathbf{x} + \iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{y} d\mathbf{x} \quad (1.30)$$

$$= \iint p(\mathbf{y}, \mathbf{x}) \ln \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}, \mathbf{x})} d\mathbf{y} d\mathbf{x} \quad (1.31)$$

$$= \iint p(\mathbf{y}, \mathbf{x}) \ln \frac{1}{p(\mathbf{x})} d\mathbf{y} d\mathbf{x} \quad (1.32)$$

$$= - \iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{x}) d\mathbf{y} d\mathbf{x} \quad (1.33)$$

$$= - \iint p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \quad (1.34)$$

$$= H[\mathbf{x}] \quad (1.35)$$

## 1.8.1 Relative entropy and mutual information

Consider some unknown distribution  $p(\mathbf{x})$ , and suppose that we have modelled this using an approximating distribution  $q(\mathbf{x})$ . Then the average **additional** amount of information (in nats) required to specify the value of  $\mathbf{x}$  as a result of using  $q(\mathbf{x})$  instead of the true distribution is given by

$$KL(q \parallel p) = - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left( - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \quad (1.36)$$

$$= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \quad (1.37)$$

This is known as the **relative entropy** or **Kullback-Leibler divergence**, or **KL divergence** (Kullback and Leibler, 1951), between the distributions  $p(\mathbf{x})$  and  $q(\mathbf{x})$ . Note that it is not a symmetrical quantity. Now introduce the concept of **convex** functions. A function  $f(x)$  is said to be convex if it has the property that every chord lies on or above the function, as shown in Figure 1.31. Any value of  $x$  in the interval from  $x = a$  to  $x = b$  can be written in the form  $a + (1 - \lambda)b$  where  $0 \leq \lambda \leq 1$ . The corresponding point on the chord is given by  $f(a) + (1 - \lambda)f(b)$ , and the corresponding value of the function is  $f(a + (1 - \lambda)b)$ . Convexity then implies

$$f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b) \quad (1.38)$$

This is equivalent to the requirement that the second derivative of the function be everywhere positive. A function is called strictly convex if the equality is satisfied only for  $\lambda = 0$  and  $\lambda = 1$ . If a function has the opposite property, namely that every chord lies on or below the function, it is called concave, with a corresponding definition for strictly concave. If a function  $f(x)$  is convex, then  $f(x)$  will be concave.

Using the technique of **proof by induction**, we can show that a convex function  $f(x)$  satisfies

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i) \quad (1.39)$$

where  $\lambda_i \geq 0$  and  $\sum_i \lambda_i = 1$ , for any set of points  $\{x_i\}$ , this is known as **Jensen's inequality**. If we interpret the  $\lambda_i$  as the probability distribution over a discrete variable  $x$  taking the values  $\{x_i\}$ , then 1.39 can be written

$$f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)] \quad (1.40)$$

where  $\mathbb{E}(\cdot)$  denotes the expectation. For continuous variables, Jensen's inequality takes the form

$$f\left(\int \mathbf{x} p(\mathbf{x})\right) \leq \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \quad (1.41)$$

We can apply Jensen's inequality to the KL-divergence to give

$$KL(p \parallel q) = -\int p(\mathbf{x}) \ln\left\{\frac{q(\mathbf{x})}{p(\mathbf{x})}\right\} d\mathbf{x} \geq -\ln \int q(\mathbf{x}) d\mathbf{x} = 0 \quad (1.42)$$

where we have used the fact that  $-\ln x$  is a convex function, together with the normalization condition  $\int q(\mathbf{x}) d\mathbf{x} = 1$ .

Suppose we try to approximate this distribution using some parametric distribution  $q(\mathbf{x}|\theta)$ . Then the expectation of KL-divergence with respect to  $p(\mathbf{x})$  can be approximated by a finite sum over these training points, so that

$$KL(p \parallel q) \simeq \sum_{n=1}^N \{-\ln q(\mathbf{x}_n|\theta) + \ln p(\mathbf{x}_n)\} \quad (1.43)$$

The second term on the right-hand side is independent of  $\theta$ , and the first term is the negative log likelihood function for  $\theta$  under the distribution  $q(\mathbf{x}|\theta)$  evaluated using the training set. Thus minimizing this KL-divergence is equivalent to maximizing the likelihood function.

If, for a joint distribution, the variables  $\mathbf{x}$  and  $\mathbf{y}$  are not independent, we can gain some idea of whether they are 'close' to being independent by considering the KL-divergence between the joint distribution and the product of the marginals, given by

$$I[\mathbf{x}, \mathbf{y}] \equiv KL(p(\mathbf{x}, \mathbf{y}) \parallel p(\mathbf{x})p(\mathbf{y})) \quad (1.44)$$

$$= -\iint p(\mathbf{x}, \mathbf{y}) \ln\left(\frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})}\right) d\mathbf{x} d\mathbf{y} \quad (1.45)$$

which is called the **mutual information** between the variables  $\mathbf{x}$  and  $\mathbf{y}$ . From the property of the KL-divergence, we see that  $I(\mathbf{x}, \mathbf{y}) \geq 0$  with equality if, and only if  $\mathbf{x}$  and  $\mathbf{y}$  are independent. Using the sum and product rules of probability, we see the relationship to conditional entropy through

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}] \quad (1.46)$$

Thus we can view the mutual information as the reduction in the uncertainty about  $\mathbf{x}$  by virtue of being told the value of  $\mathbf{y}$  (or vice versa).

## 1.9 Some basic concepts

### 1.9.1 Parametric vs non-parametric models

### 1.9.2 A simple non-parametric classifier: K-nearest neighbours

#### 1.9.2.1 Representation

$$y = f(\mathbf{x}) = \arg \min_c \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} \mathbb{I}(y_i = c) \quad (1.47)$$

where  $N_k(\mathbf{x})$  is the set of  $k$  points that are closest to point  $\mathbf{x}$ .

Usually use **k-d tree** to accelerate the process of finding  $k$  nearest points.

### 1.9.2.2 Evaluation

No training is needed.

### 1.9.2.3 Optimization

No training is needed.

## 1.9.3 Overfitting

### 1.9.4 Model selection

When we have a variety of models of different complexity (e.g., linear or logistic regression models with different degree polynomials, or KNN classifiers with different values of  $K$ ), how should we pick the right one? A natural approach is to compute the **misclassification rate** on the training set for each method.



## Chapter 2

# Probability and Statistics

One role for the distributions is to model the probability distribution  $p(x)$  of a random variable  $x$ , given a finite set of observations. This problem is known as **density estimation**. We shall assume that the data points are independent and identically distributed.

**Parameter** distributions are governed by a small number of adaptive parameters, such as the mean and variance in the case of a Gaussian for example. In a frequentist treatment, we choose specific values for the parameters by optimizing some criterion, such as the likelihood function. By contrast, in a Bayesian treatment we introduce prior distributions over the parameters and then use Bayes theorem to compute the corresponding posterior distribution given the observed data.

An important role is played by **conjugate priors**, that lead to posterior distributions having the same functional form as the prior, and that therefore lead to a greatly simplified Bayesian analytics.

**Nonparametric** density estimation methods in which the form of the distribution typically depends on the size of the data set. Parameters in such models control the model complexity rather than the form of the distribution. We cover three nonparametric methods based respectively on histograms, nearest-neighbours, and kernels.

### 2.1 Frequentists vs. Bayesians

There are two different interpretations of probability. One is called the **frequentist** interpretation. In this view, probabilities represent long run **frequencies** of events. For example, the above statement means that, if we flip the coin many times, we expect it to land heads about half the time.

The other interpretation is called the **Bayesian** interpretation of probability. In this view, probability is used to quantify our **uncertainty** about something; hence it is fundamentally related to information rather than repeated trials (Jaynes 2003). In the Bayesian view, the above statement means we believe the coin is equally likely to land heads or tails on the next toss.

One big advantage of the Bayesian interpretation is that it can be used to model our uncertainty about events that do not have long term frequencies. For example, we might want to compute the probability that the polar ice cap will melt by 2020 CE. This event will happen zero or one times, but cannot happen repeatedly. Nevertheless, we thought to be able to quantify our uncertainty about this event. To give another machine learning oriented example, we might have observed a blip on our radar screen, and want to compute the probability distribution over the location of the corresponding target (be it a bird, plane, or missile). In all these cases, the idea of repeated trials does not make sense, but the Bayesian interpretation is valid and indeed quite natural. We shall therefore adopt the Bayesian interpretation in this book. Fortunately, the basic rules of probability theory are the same, no matter which interpretation is adopted.

### 2.2 probability theory

#### 2.2.1 Concepts

The expression  $p(A)$  denotes the probability that event  $A$  is true. We require that  $0 \leq p(A) \leq 1$ , where 0 means the event definitely will not happen, and  $p(A) = 1$  means the event definitely will happen.  $p(\bar{A})$  denotes the probability of the event not  $A$ ; this is defined to be  $p(\bar{A}) = 1 - p(A)$ .

We denote a random event by defining a **random variable**  $X$ .

**Discrete random variable:**  $X$ , which can take on any value from a finite or countably infinite set. We denote the probability of the event that  $X = x$  by  $p(X = x)$ , or just  $p(x)$  for short. Here  $p(\cdot)$  is called a **probability mass function** or **pmf**. The pmfs are defined on one **state space**.  $\mathbb{I}$  denotes the binary **indicator function**.

**Continuous random variable:** the value of  $X$  is real-valued.

**prior probability distribution** In Bayesian statistical inference, a prior probability distribution, often called simply the prior, of an uncertain quantity  $p$  is the probability distribution that would express one's uncertainty about  $p$  **before** some evidence is taken into account.

**posterior probability distribution** In Bayesian statistics, the posterior probability of a random event or an uncertain proposition is the conditional probability that is assigned **after** the relevant evidence or background is taken into account.

**likelihood function** In statistics, a likelihood function (often simply the likelihood) is a function of the parameters of a statistical model. The likelihood of a set of parameter values,  $\theta$ , given outcomes  $x$ , is equal to the probability of those observed outcomes given those parameter values

## 2.2.2 Fundamental rules

In this section, we review the basic rule of probability.

### 2.2.2.1 Probability of a union of two events

Given two events,  $A$  and  $B$ , we define the probability of  $A$  or  $B$  as follows:

$$p(A \cup B) = p(A) + p(B) - p(A \cap B) \quad (2.1)$$

$$= p(A) + p(B) \quad (2.2)$$

if  $A$  and  $B$  are mutually independent

### 2.2.2.2 Joint probabilities

We define the probability of the joint event  $A$  and  $B$  as follows:

$$p(A, B) = p(A \cap B) = p(A|B)p(B) \quad (2.3)$$

This is sometimes called the **product rule**

### 2.2.2.3 Conditional probability

Define the **conditional probability** of event  $A$ , given that event  $B$  is true, as follows:

$$p(A|B) = \frac{p(A, B)}{p(B)}, \text{ if } p(B) > 0 \quad (2.4)$$

### 2.2.2.4 sum rule

$$p(X) = \sum_Y p(X, Y) \quad (2.5)$$

### 2.2.2.5 product rule

$$p(X, Y) = p(Y|X)p(X) \quad (2.6)$$

### 2.2.2.6 CDF

$$F(x) \triangleq P(X \leq x) = \begin{cases} \sum_{u \leq x} p(u) & , \text{ discrete} \\ \int_{-\infty}^x f(u) du & , \text{ continuous} \end{cases} \quad (2.7)$$

### 2.2.2.7 PMF and PDF

For discrete random variable, We denote the probability of the event that  $X = x$  by  $P(X = x)$ , or just  $p(x)$  for short. Here  $p(x)$  is called a **probability mass function** or **PMF**. A probability mass function is a function that gives the probability that a discrete random variable is exactly equal to some value<sup>4</sup>. This satisfies the properties  $0 \leq p(x) \leq 1$  and  $\sum_{x \in \mathcal{X}} p(x) = 1$ .

For continuous variable, in the equation  $F(x) = \int_{-\infty}^x f(u) du$ , the function  $f(x)$  is called a **probability density function** or **PDF**. A probability density function is a function that describes the relative likelihood for this random variable to take on a given value<sup>5</sup>. This satisfies the properties  $f(x) \geq 0$  and  $\int_{-\infty}^{\infty} f(x) dx = 1$ .

### 2.2.2.8 probability densities

probability densities

$$p(x \in (a, b)) = \int_a^b p(x) dx \quad (2.8)$$

The probability density function  $p(x)$  must satisfy the two conditions

$$\begin{cases} p(x) \geq 0 \\ \int_{-\infty}^{\infty} p(x) dx = 1 \end{cases} \quad (2.9)$$

Combinations of discrete and continuous variables.

$$p(x) = \int_p (x, y) dy \quad (2.10)$$

$$p(x, y) = p(y|x)p(x) \quad (2.11)$$

## 2.2.3 Multivariate random variables

### 2.2.3.1 Joint CDF

We denote joint CDF by  $F(x, y) \triangleq P(X \leq x \cap Y \leq y) = P(X \leq x, Y \leq y)$ .

$$F(x, y) \triangleq P(X \leq x, Y \leq y) = \begin{cases} \sum_{u \leq x, v \leq y} p(u, v) \\ \int_{-\infty}^x \int_{-\infty}^y f(u, v) du dv \end{cases} \quad (2.12)$$

**product rule:**

$$p(X, Y) = P(X|Y)P(Y) \quad (2.13)$$

**Chain rule:**

$$p(X_{1:N}) = p(X_1)p(X_2|X_1) \dots p(X_N|X_{1:N-1}) \quad (2.14)$$

<sup>4</sup> [http://en.wikipedia.org/wiki/Probability\\_mass\\_function](http://en.wikipedia.org/wiki/Probability_mass_function)

<sup>5</sup> [http://en.wikipedia.org/wiki/Probability\\_density\\_function](http://en.wikipedia.org/wiki/Probability_density_function)

### 2.2.3.2 Marginal distribution

**Marginal CDF:**

$$F_X(x) \triangleq F(x, +\infty) = \begin{cases} \sum_{x_i \leq x} P(X = x_i) = \sum_{x_i \leq x} \sum_{j=1}^{+\infty} P(X = x_i, Y = y_j) \\ \int_{-\infty}^x f_X(u) du = \int_{-\infty}^x \int_{-\infty}^{+\infty} f(u, v) du dv \end{cases} \quad (2.15)$$

$$F_Y(y) \triangleq F(+\infty, y) = \begin{cases} \sum_{y_j \leq y} p(Y = y_j) = \sum_{i=1}^{+\infty} \sum_{y_j \leq y} P(X = x_i, Y = y_j) \\ \int_{-\infty}^y f_Y(v) dv = \int_{-\infty}^{+\infty} \int_{-\infty}^y f(u, v) du dv \end{cases} \quad (2.16)$$

**Marginal PMF and PDF:**

$$\begin{cases} P(X = x_i) = \sum_{j=1}^{+\infty} P(X = x_i, Y = y_j) & , \text{ discrete} \\ f_X(x) = \int_{-\infty}^{+\infty} f(x, y) dy & , \text{ continuous} \end{cases} \quad (2.17)$$

$$\begin{cases} p(Y = y_j) = \sum_{i=1}^{+\infty} P(X = x_i, Y = y_j) & , \text{ discrete} \\ f_Y(y) = \int_{-\infty}^{+\infty} f(x, y) dx & , \text{ continuous} \end{cases} \quad (2.18)$$

### 2.2.3.3 Conditional distribution

**Conditional PMF:**

$$p(X = x_i | Y = y_j) = \frac{p(X = x_i, Y = y_j)}{p(Y = y_j)} \text{ if } p(Y) > 0 \quad (2.19)$$

The pmf  $p(X|Y)$  is called **conditional probability**.

**Conditional PDF:**

$$f_{X|Y}(x|y) = \frac{f(x, y)}{f_Y(y)} \quad (2.20)$$

## 2.2.4 Bayes rule

Bayes' theorem

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)} \quad (2.21)$$

Denominator in Bayes' theorem

$$p(X) = \sum_Y p(X|Y)p(Y) \quad (2.22)$$

**Bayesian probabilities** So far, we have viewed probabilities in terms of the frequencies of random, repeatable events, which we shall refer to as the classical or frequentist interpretation of probability. Now we turn to the more general Bayesian view, in which probabilities provide a quantification of uncertainty.

We can adopt a similar approach when making inferences about quantities such as the parameters  $\mathbf{w}$  in the polynomial curve fitting. We capture our assumptions about  $\mathbf{w}$ , before observing the data, in the form of a **prior** probability distribution  $p(\mathbf{w})$ . The effect of the observed data  $\mathcal{D} = t_1, \dots, t_N$  is expressed through the conditional probability  $p(\mathcal{D}|\mathbf{w})$ . Bayes' theorem, which takes the form

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})} \quad (2.23)$$



then allows us to evaluate the uncertainty in  $\mathbf{w}$  **after** we have observed  $\mathcal{D}$  in the form of the **posterior** probability  $p(\mathbf{w}|\mathcal{D})$ .

The quantity  $p(\mathcal{D}|\mathbf{w})$  on the right-hand side of Bayes' theorem is evaluated for the observed data set  $\mathcal{D}$  and can be viewed as a function of the parameter vector  $\mathbf{w}$ , in which case it is called the **likelihood function**. It expresses how probable the observed data set is for different settings of the parameter vector  $\mathbf{w}$ .

$$posterior \propto likelihood \times prior \quad (2.24)$$

where all of these quantities are viewed as functions of  $\mathbf{w}$ . Summing both side with respect to  $\mathbf{w}$

$$p(\mathcal{D})p(\mathbf{w}|\mathcal{D}) = p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) \quad (2.25)$$

$$\Rightarrow \sum_{\mathbf{w}} p(\mathcal{D})p(\mathbf{w}|\mathcal{D}) = \sum_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) \quad (2.26)$$

$$\Rightarrow p(\mathcal{D}) = \sum_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) \quad (2.27)$$

$$(2.28)$$

Integrating both side with respect to  $\mathbf{w}$

$$\int p(\mathcal{D})p(\mathbf{w}|\mathcal{D})d\mathbf{w} = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w} \quad (2.29)$$

$$p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w} \quad (2.30)$$

A widely used frequentist estimator is **maximum likelihood**, in which  $\mathbf{w}$  is set to the value that maximizes the likelihood function  $p(\mathcal{D}|\mathbf{w})$ . In the machine learning literature, the negative log of the likelihood function is called an **error function**.

One approach to determining the frequentist error bars is the **bootstrap**, in which multiple data sets are created as follows. Suppose our original data set consists of  $N$  data points. We can create a new data set  $X_B$  by drawing  $N$  points at random from  $X$ , with **replacement**, so that some points in  $X$  may be replicated in  $X_B$ , whereas other points in  $X$  may be absent from  $X_B$ . This process can be repeated  $L$  times to generate  $L$  data sets each of size  $N$  and each obtained by sampling from the original data set  $X$ . The statistical accuracy of parameter estimates can then be evaluated by looking at the variability of predictions between the different bootstrap data sets.

## 2.2.5 Independence and conditional independence

We say  $X$  and  $Y$  are unconditionally independent or marginally independent, denoted  $X \perp Y$ , if we can represent the joint as the product of the two marginals, i.e.,

$$X \perp Y = P(X, Y) = P(X)P(Y) \quad (2.31)$$

We say  $X$  and  $Y$  are conditionally independent (CI) given  $Z$  if the conditional joint can be written as a product of conditional marginals:

$$X \perp Y|Z = P(X, Y|Z) = P(X|Z)P(Y|Z) \quad (2.32)$$

## 2.2.6 Quantiles

Since the cdf  $F$  is a monotonically increasing function, it has an inverse; let us denote this by  $F^{-1}$ . If  $F$  is the cdf of  $X$ , then  $F^{-1}(\alpha)$  is the value of  $x_\alpha$  such that  $P(X \leq x_\alpha) = \alpha$ ; this is called the  $\alpha$  quantile of  $F$ . The value  $F^{-1}(0.5)$  is the **median** of the distribution, with half of the probability mass on the left, and half on the right. The values  $F^{-1}(0.25)$  and  $F^{-1}(0.75)$  are the lower and upper **quartiles**.

### 2.2.7 Mean and variance

The most familiar property of a distribution is its **mean**, or **expected value**, denoted by  $\mu$ . For discrete rvs, it is defined as  $\mathbb{E}[X] \triangleq \sum_{x \in \mathcal{X}} xp(x)$ , and for continuous rvs, it is defined as  $\mathbb{E}[X] \triangleq \int_{\mathcal{X}} xp(x)dx$ . If this integral is not finite, the mean is not defined (we will see some examples of this later).

The **variance** is a measure of the spread of a distribution, denoted by  $\sigma^2$ . This is defined as follows:

$$\text{var}[X] = \mathbb{E}[(X - \mu)^2] \quad (2.33)$$

$$\begin{aligned} &= \int (x - \mu)^2 p(x) dx \\ &= \int x^2 p(x) dx + \mu^2 \int p(x) dx - 2\mu \int xp(x) dx \\ &= \mathbb{E}[X^2] - \mu^2 \end{aligned} \quad (2.34)$$

from which we derive the useful result

$$\mathbb{E}[X^2] = \sigma^2 + \mu^2 \quad (2.35)$$

The **standard deviation** is defined as

$$\text{std}[X] \triangleq \sqrt{\text{var}[X]} \quad (2.36)$$

This is useful since it has the same units as  $X$  itself.

**Expectations and covariances** The average value of some function  $f(x)$  under a probability distribution  $p(x)$  is called the expectation of  $f(x)$  and will be denoted by

$$\mathbb{E}[f] = \sum_x p(x) f(x) \mathbb{E}[f] = \int p(x) f(x) dx \quad (2.37)$$

approximation

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n) \quad (2.38)$$

conditional expectation with respect to a conditional distribution

$$\mathbb{E}_x[f|y] = \sum_x p(x|y) f(x) \quad (2.39)$$

variance of  $f(x)$  is defined by

$$\text{var}[f] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] \text{var}[f] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2 \quad (2.40)$$

covariance

$$\text{cov}[x, y] = \mathbb{E}_{x,y}[\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}] \quad (2.41)$$

In the case of two vectors of random variables  $\mathbf{x}$  and  $\mathbf{y}$

$$\text{cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{x,y}[\{\mathbf{x} - \mathbb{E}[\mathbf{x}]\}\{\mathbf{y}^T - \mathbb{E}[\mathbf{y}^T]\}] \text{cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{x,y}[\mathbf{x}\mathbf{y}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}^T] \quad (2.42)$$

## 2.3 Some common discrete distributions

In this section, we review some commonly used parametric distributions defined on discrete state spaces, both finite and countably infinite.

### 2.3.1 The Bernoulli and binomial distributions

**Definition 2.1.** Now suppose we toss a coin only once. Let  $X \in \{0, 1\}$  be a binary random variable, with probability of success or heads of  $\theta$ . We say that  $X$  has a **Bernoulli distribution**. This is written as  $X \sim \text{Ber}(\theta)$ , where the pmf is defined as

$$\text{Ber}(x|\theta) \triangleq \theta^{\mathbb{I}(x=1)}(1-\theta)^{\mathbb{I}(x=0)} \quad (2.43)$$

**Definition 2.2.** Suppose we toss a coin  $n$  times. Let  $X \in \{0, 1, \dots, n\}$  be the number of heads. If the probability of heads is  $\theta$ , then we say  $X$  has a **binomial distribution**, written as  $X \sim \text{Bin}(n, \theta)$ . The pmf is given by

$$\text{Bin}(k|n, \theta) \triangleq \binom{n}{k} \theta^k (1-\theta)^{n-k} \quad (2.44)$$

### 2.3.2 The multinoulli and multinomial distributions

**Definition 2.3.** The Bernoulli distribution can be used to model the outcome of one coin tosses. To model the outcome of tossing a  $K$ -sided dice, let  $\mathbf{x} = (\mathbb{I}(x=1), \dots, \mathbb{I}(x=K)) \in \{0, 1\}^K$  be a random vector (this is called **dummy encoding** or **one-hot encoding**), then we say  $X$  has a **multinoulli distribution** (or **categorical distribution**), written as  $X \sim \text{Cat}(\theta)$ . The pmf is given by:

$$p(\mathbf{x}) \triangleq \prod_{k=1}^K \theta_k^{\mathbb{I}(x_k=1)} \quad (2.45)$$

**Definition 2.4.** Suppose we toss a  $K$ -sided dice  $n$  times. Let  $\mathbf{x} = (x_1, x_2, \dots, x_K) \in \{0, 1, \dots, n\}^K$  be a random vector, where  $x_j$  is the number of times side  $j$  of the dice occurs, then we say  $X$  has a **multinomial distribution**, written as  $X \sim \text{Mu}(n, \theta)$ . The pmf is given by

$$p(\mathbf{x}) \triangleq \binom{n}{x_1 \dots x_K} \prod_{k=1}^K \theta_k^{x_k} \quad (2.46)$$

where  $\binom{n}{x_1 \dots x_K} \triangleq \frac{n!}{x_1! x_2! \dots x_K!}$

Bernoulli distribution is just a special case of a Binomial distribution with  $n = 1$ , and so is multinoulli distribution as to multinomial distribution. See Table 2.1 for a summary.

Table 2.1: Summary of the multinomial and related distributions.

Name	K	n	X
Bernoulli	1	1	$x \in \{0, 1\}$
Binomial	1	-	$\mathbf{x} \in \{0, 1, \dots, n\}$
Multinoulli	-	1	$\mathbf{x} \in \{0, 1\}^K, \sum_{k=1}^K x_k = 1$
Multinomial	-	-	$\mathbf{x} \in \{0, 1, \dots, n\}^K, \sum_{k=1}^K x_k = n$

### 2.3.3 The Poisson distribution

**Definition 2.5.** We say that  $X \in \{0, 1, 2, \dots\}$  has a **Poisson distribution** with parameter  $\lambda > 0$ , written as  $X \sim \text{Poi}(\lambda)$ , if its pmf is

$$p(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!} \quad (2.47)$$

The first term is just the normalization constant, required to ensure the distribution sums to 1.

The Poisson distribution is often used as a model for counts of rare events like radioactive decay and traffic accidents.

Table 2.2: Summary of Bernoulli, binomial multinoulli and multinomial distributions.

Name	Written as	X	$p(x)$ (or $p(\mathbf{x})$ )	$\mathbb{E}[X]$	$\text{var}[X]$
Bernoulli	$X \sim \text{Ber}(\theta)$	$x \in \{0, 1\}$	$\theta^{\mathbb{I}(x=1)}(1-\theta)^{\mathbb{I}(x=0)}$	$\theta$	$\theta(1-\theta)$
Binomial	$X \sim \text{Bin}(n, \theta)$	$x \in \{0, 1, \dots, n\}$	$\binom{n}{k} \theta^k (1-\theta)^{n-k}$	$n\theta$	$n\theta(1-\theta)$
Multinoulli	$X \sim \text{Cat}(\theta)$	$\mathbf{x} \in \{0, 1\}^K, \sum_{k=1}^K x_k = 1$	$\prod_{k=1}^K \theta_j^{\mathbb{I}(x_j=1)}$	-	-
Multinomial	$X \sim \text{Mu}(n, \theta)$	$\mathbf{x} \in \{0, 1, \dots, n\}^K, \sum_{k=1}^K x_k = n$	$\binom{n}{x_1 \dots x_K} \prod_{k=1}^K \theta_j^{x_j}$	-	-
Poisson	$X \sim \text{Poi}(\lambda)$	$x \in \{0, 1, 2, \dots\}$	$e^{-\lambda} \frac{\lambda^x}{x!}$	$\lambda$	$\lambda$

### 2.3.4 The empirical distribution

The **empirical distribution function**<sup>6</sup>, or **empirical cdf**, is the cumulative distribution function associated with the empirical measure of the sample. Let  $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$  be a sample set, it is defined as

$$F_n(x) \triangleq \frac{1}{N} \sum_{i=1}^N \mathbb{I}(x_i \leq x) \quad (2.48)$$

## 2.4 Some common continuous distributions

In this section we present some commonly used univariate (one-dimensional) continuous probability distributions.

### 2.4.1 Gaussian (normal) distribution

Table 2.3: Summary of Gaussian distribution.

Written as	$f(x)$	$\mathbb{E}[X]$	mode	$\text{var}[X]$
$X \sim \mathcal{N}(\mu, \sigma^2)$	$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$	$\mu$	$\mu$	$\sigma^2$

If  $X \sim \mathcal{N}(0, 1)$ , we say  $X$  follows a **standard normal** distribution.

The Gaussian distribution is the most widely used distribution in statistics. There are several reasons for this.

1. First, it has two parameters which are easy to interpret, and which capture some of the most basic properties of a distribution, namely its mean and variance.

<sup>6</sup> [http://en.wikipedia.org/wiki/Empirical\\_distribution\\_function](http://en.wikipedia.org/wiki/Empirical_distribution_function)

2. Second, the central limit theorem (Section TODO) tells us that sums of independent random variables have an approximately Gaussian distribution, making it a good choice for modeling residual errors or noise.
3. Third, the Gaussian distribution makes the least number of assumptions (has maximum entropy), subject to the constraint of having a specified mean and variance, as we show in Section TODO; this makes it a good default choice in many cases.
4. Finally, it has a simple mathematical form, which results in easy to implement, but often highly effective, methods, as we will see.

See (Jaynes 2003, ch 7) for a more extensive discussion of why Gaussians are so widely used. More about Gaussian distribution is discussed in 2.5.

## 2.4.2 Student's t-distribution

Table 2.4: Summary of Student's t-distribution.

Written as	$f(x)$	$\mathbb{E}[X]$	mode	$\text{var}[X]$
$X \sim \mathcal{T}(\mu, \sigma^2, \nu)$	$\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left[ 1 + \frac{1}{\nu} \left( \frac{x-\mu}{\sigma} \right)^2 \right]^{-\frac{\nu+1}{2}}$	$\mu$	$\mu$	$\frac{\nu\sigma^2}{\nu-2}$

where  $\Gamma(x)$  is the gamma function:

$$\Gamma(x) \triangleq \int_0^\infty t^{x-1} e^{-t} dt \quad (2.49)$$

$\mu$  is the mean,  $\sigma^2 > 0$  is the scale parameter, and  $\nu > 0$  is called the **degrees of freedom**. See Figure 2.1 for some plots.

The variance is only defined if  $\nu > 2$ . The mean is only defined if  $\nu > 1$ .

As an illustration of the robustness of the Student distribution, consider Figure 2.2. We see that the Gaussian is affected a lot, whereas the Student distribution hardly changes. This is because the Student has heavier tails, at least for small  $\nu$  (see Figure 2.1).

If  $\nu = 1$ , this distribution is known as the **Cauchy** or **Lorentz** distribution. This is notable for having such heavy tails that the integral that defines the mean does not converge.

To ensure finite variance, we require  $\nu > 2$ . It is common to use  $\nu = 4$ , which gives good performance in a range of problems (Lange et al. 1989). For  $\nu \gg 5$ , the Student distribution rapidly approaches a Gaussian distribution and loses its robustness properties.

## 2.4.3 The Laplace distribution

Table 2.5: Summary of Laplace distribution.

Written as	$f(x)$	$\mathbb{E}[X]$	mode	$\text{var}[X]$
$X \sim \text{Lap}(\mu, b)$	$\frac{1}{2b} \exp\left(-\frac{ x-\mu }{b}\right)$	$\mu$	$\mu$	$2b^2$

Here  $\mu$  is a location parameter and  $b > 0$  is a scale parameter. See Figure 2.1 for a plot.

Its robustness to outliers is illustrated in Figure 2.2. It also puts more probability density at 0 than the Gaussian. This property is a useful way to encourage sparsity in a model, as we will see in Section TODO.



Fig. 2.1: (a) The pdfs for a  $\mathcal{N}(0, 1)$ ,  $\mathcal{T}(0, 1, 1)$  and  $Lap(0, 1/\sqrt{2})$ . The mean is 0 and the variance is 1 for both the Gaussian and Laplace. The mean and variance of the Student is undefined when  $\nu = 1$ . (b) Log of these pdfs. Note that the Student distribution is not log-concave for any parameter value, unlike the Laplace distribution, which is always log-concave (and log-convex...) Nevertheless, both are unimodal.

## 2.4.4 The gamma distribution

Table 2.6: Summary of gamma distribution

Written as	$X$	$f(x)$	$\mathbb{E}[X]$	mode	$\text{var}[X]$
$X \sim \text{Ga}(a, b)$	$x \in \mathbb{R}^+$	$\frac{b^a}{\Gamma(a)} x^{a-1} e^{-xb}$	$\frac{a}{b}$	$\frac{a-1}{b}$	$\frac{a}{b^2}$

Here  $a > 0$  is called the shape parameter and  $b > 0$  is called the rate parameter. See Figure 2.3 for some plots.



Fig. 2.2: Illustration of the effect of outliers on fitting Gaussian, Student and Laplace distributions. (a) No outliers (the Gaussian and Student curves are on top of each other). (b) With outliers. We see that the Gaussian is more affected by outliers than the Student and Laplace distributions.

## 2.4.5 The beta distribution

Table 2.7: Summary of Beta distribution

Name	Written as	$X$	$f(x)$	$\mathbb{E}[X]$	mode	$\text{var}[X]$
Beta distribution	$X \sim \text{Beta}(a, b)$	$x \in [0, 1]$	$\frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1}$	$\frac{a}{a+b}$	$\frac{a-1}{a+b-2}$	$\frac{ab}{(a+b)^2(a+b+1)}$

Here  $B(a, b)$  is the beta function,

$$B(a, b) \triangleq \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \quad (2.50)$$

See Figure 2.4 for plots of some beta distributions. We require  $a, b > 0$  to ensure the distribution is integrable (i.e., to ensure  $B(a, b)$  exists). If  $a = b = 1$ , we get the uniform distribution. If  $a$  and  $b$  are both less than 1, we get a bimodal distribution with spikes at 0 and 1; if  $a$  and  $b$  are both greater than 1, the distribution is unimodal.



Fig. 2.3: Some  $\text{Ga}(a, b = 1)$  distributions. If  $a \leq 1$ , the mode is at 0, otherwise it is  $> 0$ . As we increase the rate  $b$ , we reduce the horizontal scale, thus squeezing everything leftwards and upwards. (b) An empirical pdf of some rainfall data, with a fitted Gamma distribution superimposed.



Fig. 2.4: Some beta distributions.



## 2.4.6 Pareto distribution

Table 2.8: Summary of Pareto distribution

Name	Written as	$X$	$f(x)$	$\mathbb{E}[X]$	mode	$\text{var}[X]$
Pareto distribution	$X \sim \text{Pareto}(k, m)$	$x \geq m$	$km^k x^{-(k+1)} \mathbb{I}(x \geq m)$	$\frac{km}{k-1}$ if $k > 1$	$m$	$\frac{m^2 k}{(k-1)^2(k-2)}$ if $k > 2$

The **Pareto distribution** is used to model the distribution of quantities that exhibit **long tails**, also called **heavy tails**.

As  $k \rightarrow \infty$ , the distribution approaches  $\delta(x - m)$ . See Figure 2.5(a) for some plots. If we plot the distribution on a log-log scale, it forms a straight line, of the form  $\log p(x) = a \log x + c$  for some constants  $a$  and  $c$ . See Figure 2.5(b) for an illustration (this is known as a **power law**).

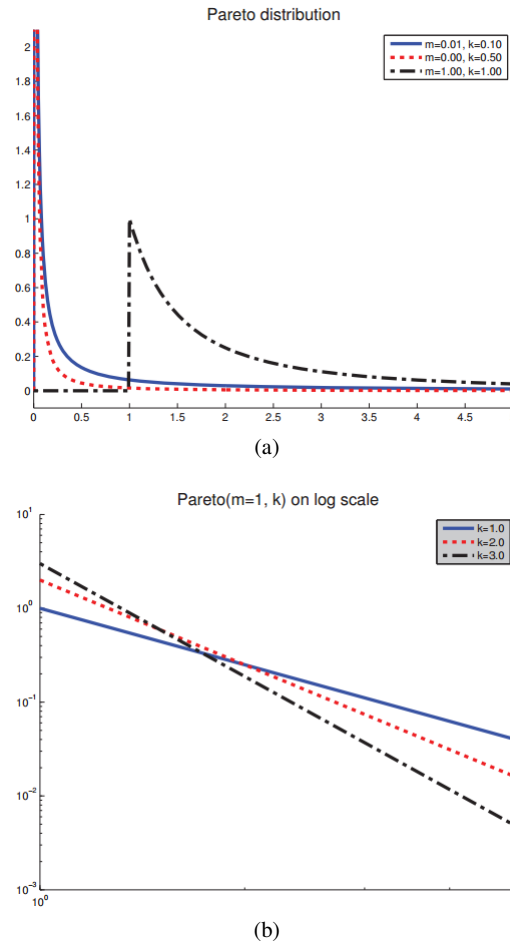


Fig. 2.5: (a) The Pareto distribution  $\text{Pareto}(x|m, k)$  for  $m = 1$ . (b) The pdf on a log-log scale.

## 2.5 The Gaussian Distribution

The Gaussian distributions, also known as the normal distributions, is a widely used model for the distribution of continuous variables. In the case of single variable  $x$ , the Gaussian is in the form

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{(2\sigma)^2}(x-\mu)^2\right\} \quad (2.51)$$

where  $\mu$  is the mean and  $\sigma^2$  is the variance. For a  $D$ -dimensional vector  $\mathbf{x}$ , the multivariate Gaussian distribution takes the form

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\} \quad (2.52)$$

where  $\boldsymbol{\mu}$  is a  $D$ -dimensional mean vector,  $\boldsymbol{\Sigma}$  is a  $D \times D$  covariance matrix, and  $|\boldsymbol{\Sigma}|$  denotes the determinant of  $\boldsymbol{\Sigma}$ .

The Gaussian distribution arises in many different contexts and can be motivated from a variety of different perspectives, such as maximizing the entropy and the sum of multiple random variables. The **central limit theorem** tells us the latter.

### 2.5.1 Conditional Gaussian distributions

### 2.5.2 Marginal Gaussian distributions

### 2.5.3 Baye's theorem for Gaussian variables

## 2.6 The exponential family

Before defining the exponential family, we mention several reasons why it is important:

- It can be shown that, under certain regularity conditions, the exponential family is the only family of distributions with finite-sized sufficient statistics, meaning that we can compress the data into a fixed-sized summary without loss of information. This is particularly useful for online learning, as we will see later.
- The exponential family is the only family of distributions for which conjugate priors exist, which simplifies the computation of the posterior (see Section 2.6.8).
- The exponential family can be shown to be the family of distributions that makes the least set of assumptions subject to some user-chosen constraints (see Section 2.6.9).
- The exponential family is at the core of generalized linear models, as discussed in Section ??.
- The exponential family is at the core of variational inference, as discussed in Section TODO.

### 2.6.1 Definition

A pdf or pmf  $p(\mathbf{x}|\boldsymbol{\theta})$ , for  $\mathbf{x} \in \mathbb{R}^m$  and  $\boldsymbol{\theta} \in \mathbb{R}^D$ , is said to be in the **exponential family** if it is of the form

$$p(\mathbf{x}|\boldsymbol{\theta}) = h(\mathbf{x})g(\boldsymbol{\theta}) \exp\{\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x})\} \quad (2.53)$$

$$= \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{x}) \exp[\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x})] \quad (2.54)$$

$$= h(\mathbf{x}) \exp[\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) - A(\boldsymbol{\theta})] \quad (2.55)$$

where

$$Z(\theta) = \int h(x) \exp[\theta^T \phi(x)] dx \quad (2.56)$$

$$A(\theta) = \log Z(\theta) \quad (2.57)$$

Here  $\theta$  are called the **natural parameters** or **canonical parameters**,  $\phi(x) \in \mathbb{R}^D$  is called a vector of **sufficient statistics**,  $Z(\theta)$  is called the **partition function**,  $A(\theta)$  is called the **log partition function** or **cumulant function**, and  $h(x)$  is the a scaling constant, often 1. If  $\phi(x) = x$ , we say it is a **natural exponential family**.

Equation 2.55 can be generalized by writing

$$p(x|\theta) = h(x) \exp[\eta(\theta)^T \phi(x) - A(\eta(\theta))] \quad (2.58)$$

where  $\eta$  is a function that maps the parameters  $\theta$  to the canonical parameters  $\eta = \eta(\theta)$ . If  $\dim(\theta) < \dim(\eta(\theta))$ , it is called a **curved exponential family**, which means we have more sufficient statistics than parameters. If  $\eta(\theta) = \theta$ , the model is said to be in **canonical form**. We will assume models are in canonical form unless we state otherwise.

## 2.6.2 Maximum likelihood and sufficient statistics

## 2.6.3 Conjugate priors

## 2.6.4 Noninformative priors

In many cases, however, we may have little idea of what form the distribution should take. We may then seek a form of prior distribution, called a **noninformative prior**, which is intended to have little influence on the posterior distribution as possible (Jeffries, 1946; Box and Tiao, 1973; Bernardo and Smith, 1994).

Suppose we have a distribution  $p(x|\lambda)$  governed by a parameter  $\lambda$ , we might be tempted to propose a prior distribution  $p(\lambda) = \text{const}$  as a suitable prior. Improper priors are ones in which the domain  $\lambda$  is unbounded and the prior distribution cannot be correctly normalized because the integral over  $\lambda$  diverges. A second difficulty arises from the transformation behaviour of a probability density under a nonlinear change of variables. For example

$$p_\eta(\eta) = p_\lambda(\lambda) \left| \frac{d\lambda}{d\eta} \right| = p_\lambda(\eta^2) \eta \propto \eta \quad (2.59)$$

so the density over  $\eta$  will not be constant.

Here are two simple examples of noninformative priors (Berger, 1985). First of all, if a density takes the form

$$p(x|\mu) = f(x - \mu) \quad (2.60)$$

then the parameter  $\mu$  is known as a **location parameter**. This family of densities exhibits **translation invariance** because if we shift  $x$  by a constant to give  $\hat{x} = x + c$ , then

$$p(\hat{x}|\hat{\mu}) = f(\hat{x} - \hat{\mu}) \quad (2.61)$$

where  $\hat{\mu} = \mu + c$ .

A second example,

$$p(x|\sigma) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right) \quad (2.62)$$

where  $\sigma > 0$ . Note that this will be a normalized density provided  $f(x)$  is correctly normalized. The parameter  $\sigma$  is known as a **scale parameter**, and the density exhibits **scale invariance**.

## 2.6.5 Examples

### 2.6.5.1 Bernoulli

The Bernoulli for  $x \in \{0, 1\}$  can be written in exponential family form as follows:

$$\begin{aligned}\text{Ber}(x|\mu) &= \mu^x(1-\mu)^{1-x} \\ &= \exp[x \log \mu + (1-x) \log(1-\mu)]\end{aligned}\tag{2.63}$$

where  $\phi(x) = (\mathbb{I}(x=0), \mathbb{I}(x=1))$  and  $\theta = (\log \mu, \log(1-\mu))$ .

However, this representation is **over-complete** since  $\mathbf{1}^T \phi(x) = \mathbb{I}(x=0) + \mathbb{I}(x=1) = 1$ . Consequently  $\theta$  is not uniquely identifiable. It is common to require that the representation be **minimal**, which means there is a unique  $\theta$  associated with the distribution. In this case, we can just define

$$\begin{aligned}\text{Ber}(x|\mu) &= (1-\mu) \exp\left(x \log \frac{\mu}{1-\mu}\right) \\ \text{where } \phi(x) &= x, \theta = \log \frac{\mu}{1-\mu}, Z = \frac{1}{1-\mu}\end{aligned}\tag{2.64}$$

We can recover the mean parameter  $\mu$  from the canonical parameter using

$$\mu = \text{sigm}(\theta) = \frac{1}{1 + e^{-\theta}}\tag{2.65}$$

### 2.6.5.2 Multinoulli

We can represent the multinoulli as a minimal exponential family as follows:

$$\begin{aligned}\text{Cat}(\mathbf{x}|\boldsymbol{\mu}) &= \prod_{k=1}^K \exp\left(\sum_{k=1}^K x_k \log \mu_k\right) \\ &= \exp\left[\sum_{k=1}^{K-1} x_k \log \mu_k + \left(1 - \sum_{k=1}^{K-1} x_k\right) \log\left(1 - \sum_{k=1}^{K-1} \mu_k\right)\right] \\ &= \exp\left[\sum_{k=1}^{K-1} x_k \log \frac{\mu_k}{1 - \sum_{k=1}^{K-1} \mu_k} + \log\left(1 - \sum_{k=1}^{K-1} \mu_k\right)\right] \\ &= \exp\left[\sum_{k=1}^{K-1} x_k \log \frac{\mu_k}{\mu_K} + \log \mu_K\right], \text{ where } \mu_K \triangleq 1 - \sum_{k=1}^{K-1} \mu_k\end{aligned}$$

We can write this in exponential family form as follows:

$$\text{Cat}(\mathbf{x}|\boldsymbol{\mu}) = \exp[\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) - A(\boldsymbol{\theta})]\tag{2.66}$$

$$\boldsymbol{\theta} \triangleq \left(\log \frac{\mu_1}{\mu_K}, \dots, \log \frac{\mu_{K-1}}{\mu_K}\right)\tag{2.67}$$

$$\boldsymbol{\phi}(\mathbf{x}) \triangleq (x_1, \dots, x_{K-1})\tag{2.68}$$

We can recover the mean parameters from the canonical parameters using

$$\mu_k = \frac{e^{\theta_k}}{1 + \sum_{j=1}^{K-1} e^{\theta_j}} \quad (2.69)$$

$$\mu_K = 1 - \frac{\sum_{j=1}^{K-1} e^{\theta_j}}{1 + \sum_{j=1}^{K-1} e^{\theta_j}} = \frac{1}{1 + \sum_{j=1}^{K-1} e^{\theta_j}} \quad (2.70)$$

and hence

$$A(\boldsymbol{\theta}) = -\log \mu_K = \log\left(1 + \sum_{j=1}^{K-1} e^{\theta_j}\right) \quad (2.71)$$

### 2.6.5.3 Univariate Gaussian

The univariate Gaussian can be written in exponential family form as follows:

$$\begin{aligned} \mathcal{N}(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(x-\mu)^2\right] \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}\mu^2\right] \\ &= \frac{1}{Z(\boldsymbol{\theta})} \exp[\boldsymbol{\theta}^T \boldsymbol{\phi}(x)] \end{aligned} \quad (2.72)$$

where

$$\boldsymbol{\theta} = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right) \quad (2.73)$$

$$\boldsymbol{\phi}(x) = (x, x^2) \quad (2.74)$$

$$Z(\boldsymbol{\theta}) = \sqrt{2\pi}\sigma \exp\left(\frac{\mu^2}{2\sigma^2}\right) \quad (2.75)$$

### 2.6.5.4 Non-examples

Not all distributions of interest belong to the exponential family. For example, the uniform distribution,  $X \sim U(a, b)$ , does not, since the support of the distribution depends on the parameters. Also, the Student T distribution (Section TODO) does not belong, since it does not have the required form.

## 2.6.6 Log partition function

An important property of the exponential family is that derivatives of the log partition function can be used to generate **cumulants** of the sufficient statistics.<sup>7</sup> For this reason,  $A(\boldsymbol{\theta})$  is sometimes called a **cumulant function**. We will prove this for a 1-parameter distribution; this can be generalized to a  $K$ -parameter distribution in a straightforward way. For the first derivative we have

For the second derivative we have

---

<sup>7</sup> The first and second cumulants of a distribution are its mean  $\mathbb{E}[X]$  and variance  $\text{var}[X]$ , whereas the first and second moments are its mean  $\mathbb{E}[X]$  and  $\mathbb{E}[X^2]$ .

$$\begin{aligned}
\frac{dA}{d\theta} &= \frac{d}{d\theta} \left\{ \log \int \exp[\theta \phi(x)] h(x) dx \right\} \\
&= \frac{\frac{d}{d\theta} \int \exp[\theta \phi(x)] h(x) dx}{\int \exp[\theta \phi(x)] h(x) dx} \\
&= \frac{\int \phi(x) \exp[\theta \phi(x)] h(x) dx}{\exp(A(\theta))} \\
&= \int \phi(x) \exp[\theta \phi(x) - A(\theta)] h(x) dx \\
&= \int \phi(x) p(x) dx = \mathbb{E}[\phi(x)]
\end{aligned} \tag{2.76}$$

For the second derivative we have

$$\begin{aligned}
\frac{d^2 A}{d\theta^2} &= \int \phi(x) \exp[\theta \phi(x) - A(\theta)] h(x) [\phi(x) - A'(\theta)] dx \\
&= \int \phi(x) p(x) [\phi(x) - A'(\theta)] dx \\
&= \int \phi^2(x) p(x) dx - A'(\theta) \int \phi(x) p(x) dx \\
&= \mathbb{E}[\phi^2(x)] - \mathbb{E}[\phi(x)]^2 = \text{var}[\phi(x)]
\end{aligned} \tag{2.77}$$

In the multivariate case, we have that

$$\frac{\partial^2 A}{\partial \theta_i \partial \theta_j} = \mathbb{E}[\phi_i(x) \phi_j(x)] - \mathbb{E}[\phi_i(x)] \mathbb{E}[\phi_j(x)] \tag{2.78}$$

and hence

$$\nabla^2 A(\theta) = \text{cov}[\phi(x)] \tag{2.79}$$

Since the covariance is positive definite, we see that  $A(\theta)$  is a convex function (see Section B.1).

### 2.6.7 MLE for the exponential family

The likelihood of an exponential family model has the form

$$p(\mathcal{D}|\theta) = \left[ \prod_{i=1}^N h(x_i) \right] g(\theta)^N \exp \left[ \theta^T \left( \sum_{i=1}^N \phi(x_i) \right) \right] \tag{2.80}$$

We see that the sufficient statistics are  $N$  and

$$\phi(\mathcal{D}) = \sum_{i=1}^N \phi(x_i) = \left( \sum_{i=1}^N \phi_1(x_i), \dots, \sum_{i=1}^N \phi_K(x_i) \right) \tag{2.81}$$

The **Pitman-Koopman-Darmois theorem** states that, under certain regularity conditions, the exponential family is the only family of distributions with finite sufficient statistics. (Here, finite means of a size independent of the size of the data set.)

One of the conditions required in this theorem is that the support of the distribution not be dependent on the parameter.

## 2.6.8 Bayes for the exponential family

TODO

### 2.6.8.1 Likelihood

## 2.6.9 Maximum entropy derivation of the exponential family \*

## 2.7 Nonparametric Methods

### 2.7.1 Kernel density estimators

### 2.7.2 Nearest-neighbour methods

## 2.8 Joint probability distributions

Given a **multivariate random variable** or **random vector** <sup>8</sup>  $X \in \mathbb{R}^D$ , the **joint probability distribution** <sup>9</sup> is a probability distribution that gives the probability that each of  $X_1, X_2, \dots, X_D$  falls in any particular range or discrete set of values specified for that variable. In the case of only two random variables, this is called a **bivariate distribution**, but the concept generalizes to any number of random variables, giving a **multivariate distribution**.

The joint probability distribution can be expressed either in terms of a **joint cumulative distribution function** or in terms of a **joint probability density function** (in the case of continuous variables) or **joint probability mass function** (in the case of discrete variables).

### 2.8.1 Covariance and correlation

**Definition 2.6.** The **covariance** between two rvs  $X$  and  $Y$  measures the degree to which  $X$  and  $Y$  are (linearly) related. Covariance is defined as

$$\begin{aligned} \text{cov}[X, Y] &\triangleq \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \end{aligned} \quad (2.82)$$

**Definition 2.7.** If  $X$  is a  $D$ -dimensional random vector, its **covariance matrix** is defined to be the following symmetric, positive definite matrix:

$$\text{cov}[\mathbf{x}] \triangleq \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \quad (2.83)$$

$$= \begin{pmatrix} \text{var}[\mathbf{x}_1] & \text{Cov}[\mathbf{x}_1, \mathbf{x}_2] & \cdots & \text{Cov}[\mathbf{x}_1, \mathbf{x}_D] \\ \text{Cov}[\mathbf{x}_2, \mathbf{x}_1] & \text{var}[\mathbf{x}_2] & \cdots & \text{Cov}[\mathbf{x}_2, \mathbf{x}_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[\mathbf{x}_D, \mathbf{x}_1] & \text{Cov}[\mathbf{x}_D, \mathbf{x}_2] & \cdots & \text{var}[\mathbf{x}_D] \end{pmatrix} \quad (2.84)$$

For

**Definition 2.8.** The (Pearson) **correlation coefficient** between  $X$  and  $Y$  is defined as

<sup>8</sup> [http://en.wikipedia.org/wiki/Multivariate\\_random\\_variable](http://en.wikipedia.org/wiki/Multivariate_random_variable)

<sup>9</sup> [http://en.wikipedia.org/wiki/Joint\\_probability\\_distribution](http://en.wikipedia.org/wiki/Joint_probability_distribution)

$$\text{corr}[X, Y] \triangleq \frac{\text{Cov}[X, Y]}{\sqrt{\text{var}[X], \text{var}[Y]}} \quad (2.85)$$

A **correlation matrix** has the form

$$\mathbf{R} \triangleq \begin{pmatrix} \text{corr}[X_1, X_1] & \text{corr}[X_1, X_2] & \cdots & \text{corr}[X_1, X_D] \\ \text{corr}[X_2, X_1] & \text{corr}[X_2, X_2] & \cdots & \text{corr}[X_2, X_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{corr}[X_D, X_1] & \text{corr}[X_D, X_2] & \cdots & \text{corr}[X_D, X_D] \end{pmatrix} \quad (2.86)$$

The correlation coefficient can be viewed as a degree of linearity between  $X$  and  $Y$ , see Figure 2.6.

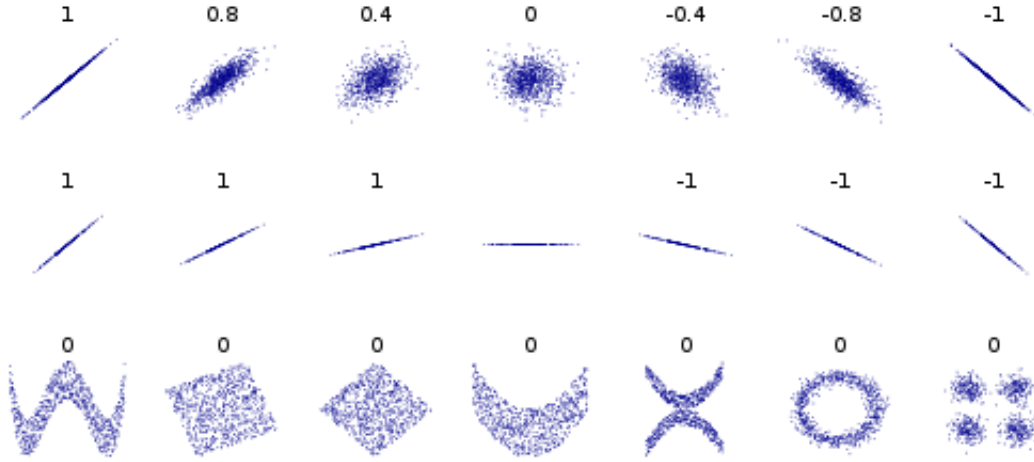


Fig. 2.6: Several sets of  $(x, y)$  points, with the Pearson correlation coefficient of  $x$  and  $y$  for each set. Note that the correlation reflects the noisiness and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many aspects of nonlinear relationships (bottom). N.B.: the figure in the center has a slope of 0 but in that case the correlation coefficient is undefined because the variance of  $Y$  is zero. Source: <http://en.wikipedia.org/wiki/Correlation>

**Uncorrelated does not imply independent.** For example, let  $X \sim U(-1, 1)$  and  $Y = X^2$ . Clearly  $Y$  is dependent on  $X$  (in fact,  $Y$  is uniquely determined by  $X$ ), yet one can show that  $\text{corr}[X, Y] = 0$ . Some striking examples of this fact are shown in Figure 2.6. This shows several data sets where there is clear dependence between  $X$  and  $Y$ , and yet the correlation coefficient is 0. A more general measure of dependence between random variables is mutual information, see Section TODO.

## 2.8.2 Multivariate Gaussian distribution

The **multivariate Gaussian** or **multivariate normal** (MVN) is the most widely used joint probability density function for continuous variables. We discuss MVNs in detail in Chapter 4; here we just give some definitions and plots.

The pdf of the MVN in  $D$  dimensions is defined by the following:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (2.87)$$

where  $\boldsymbol{\mu} = \mathbb{E}[X] \in \mathbb{R}^D$  is the mean vector, and  $\boldsymbol{\Sigma} = \text{Cov}[X]$  is the  $D \times D$  covariance matrix. The normalization constant  $(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}$  just ensures that the pdf integrates to 1.



Figure 2.7 plots some MVN densities in 2d for three different kinds of covariance matrices. A full covariance matrix has  $D(D+1)/2$  parameters (we divide by 2 since  $\Sigma$  is symmetric). A diagonal covariance matrix has  $D$  parameters, and has 0s in the off-diagonal terms. A spherical or isotropic covariance,  $\Sigma = \sigma^2 I_D$ , has one free parameter.

### 2.8.3 Multivariate Student's t-distribution

A more robust alternative to the MVN is the multivariate Student's t-distribution, whose pdf is given by

$$\begin{aligned} \mathcal{T}(x|\mu, \Sigma, \nu) &\triangleq \frac{\Gamma(\frac{\nu+D}{2})}{\Gamma(\frac{\nu}{2})} \frac{|\Sigma|^{-\frac{1}{2}}}{(\nu\pi)^{\frac{D}{2}}} \left[ 1 + \frac{1}{\nu} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]^{-\frac{\nu+D}{2}} \end{aligned} \quad (2.88)$$

$$= \frac{\Gamma(\frac{\nu+D}{2})}{\Gamma(\frac{\nu}{2})} \frac{|\Sigma|^{-\frac{1}{2}}}{(\nu\pi)^{\frac{D}{2}}} \left[ 1 + (x - \mu)^T \mathbf{V}^{-1} (x - \mu) \right]^{-\frac{\nu+D}{2}} \quad (2.89)$$

where  $\Sigma$  is called the scale matrix (since it is not exactly the covariance matrix) and  $\mathbf{V} = \nu\Sigma$ . This has fatter tails than a Gaussian. The smaller  $\nu$  is, the fatter the tails. As  $\nu \rightarrow \infty$ , the distribution tends towards a Gaussian. The distribution has the following properties

$$\text{mean} = \mu, \text{ mode} = \mu, \text{ Cov} = \frac{\nu}{\nu-2} \Sigma \quad (2.90)$$

### 2.8.4 Dirichlet distribution

A multivariate generalization of the beta distribution is the **Dirichlet distribution**, which has support over the probability simplex, defined by

$$S_K = \left\{ \mathbf{x} : 0 \leq x_k \leq 1, \sum_{k=1}^K x_k = 1 \right\} \quad (2.91)$$

The pdf is defined as follows:

$$\text{Dir}(\mathbf{x}|\alpha) \triangleq \frac{1}{B(\alpha)} \prod_{k=1}^K x_k^{\alpha_k-1} \mathbb{I}(\mathbf{x} \in S_K) \quad (2.92)$$

where  $B(\alpha_1, \alpha_2, \dots, \alpha_K)$  is the natural generalization of the beta function to  $K$  variables:

$$B(\alpha) \triangleq \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\alpha_0)} \text{ where } \alpha_0 \triangleq \sum_{k=1}^K \alpha_k \quad (2.93)$$

Figure 2.8 shows some plots of the Dirichlet when  $K = 3$ , and Figure 2.9 for some sampled probability vectors. We see that  $\alpha_0$  controls the strength of the distribution (how peaked it is), and the control where the peak occurs. For example,  $\text{Dir}(1, 1, 1)$  is a uniform distribution,  $\text{Dir}(2, 2, 2)$  is a broad distribution centered at  $(1/3, 1/3, 1/3)$ , and  $\text{Dir}(20, 20, 20)$  is a narrow distribution centered at  $(1/3, 1/3, 1/3)$ . If  $\alpha_k < 1$  for all  $k$ , we get spikes at the corner of the simplex.

For future reference, the distribution has these properties

$$\mathbb{E}(x_k) = \frac{\alpha_k}{\alpha_0}, \text{ mode}[x_k] = \frac{\alpha_k - 1}{\alpha_0 - K}, \text{ var}[x_k] = \frac{\alpha_k(\alpha_0 - \alpha_k)}{\alpha_0^2(\alpha_0 + 1)} \quad (2.94)$$

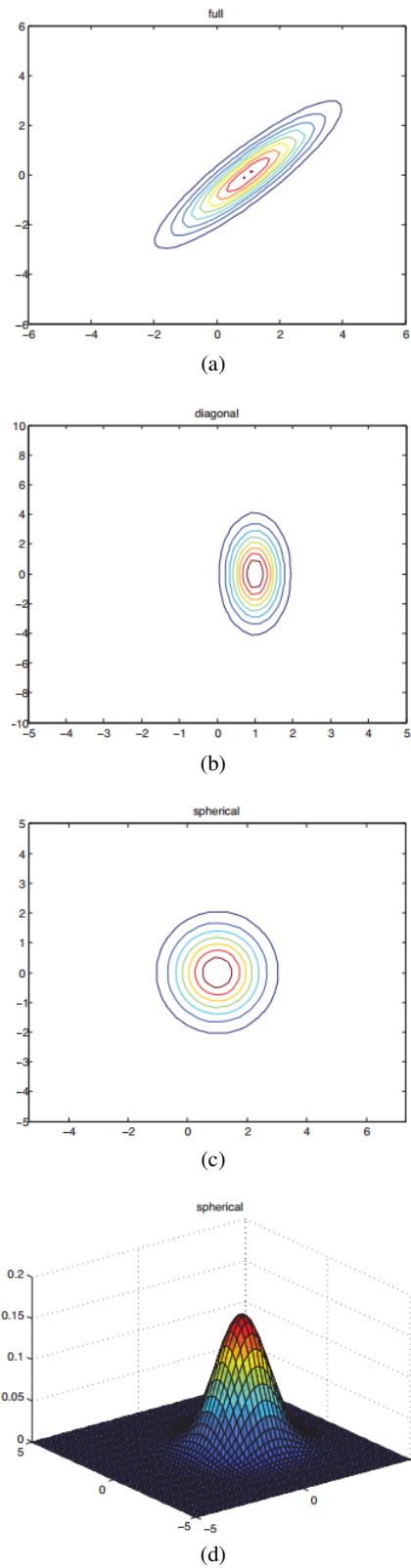


Fig. 2.7: We show the level sets for 2d Gaussians. (a) A full covariance matrix has elliptical contours. (b) A diagonal covariance matrix is an axis aligned ellipse. (c) A spherical covariance matrix has a circular shape. (d) Surface plot for the spherical Gaussian in (c).



Fig. 2.8: (a) The Dirichlet distribution when  $K = 3$  defines a distribution over the simplex, which can be represented by the triangular surface. Points on this surface satisfy  $0 \leq \theta_k \leq 1$  and  $\sum_{k=1}^K \theta_k = 1$ . (b) Plot of the Dirichlet density when  $\alpha = (2, 2, 2)$ . (c)  $\alpha = (20, 2, 2)$ .



(a)  $\alpha = (0.1, \dots, 0.1)$ . This results in very sparse distributions, with many 0s.



(b)  $\alpha = (1, \dots, 1)$ . This results in more uniform (and dense) distributions.

Fig. 2.9: Samples from a 5-dimensional symmetric Dirichlet distribution for different parameter values.

## 2.9 Transformations of random variables

If  $\mathbf{x} \sim P()$  is some random variable, and  $\mathbf{y} = f(\mathbf{x})$ , what is the distribution of  $\mathbf{Y}$ ? This is the question we address in this section.

### 2.9.1 Linear transformations

Suppose  $g()$  is a linear function:

$$g(\mathbf{x}) = A\mathbf{x} + b \quad (2.95)$$

First, for the mean, we have

$$\mathbb{E}[\mathbf{y}] = \mathbb{E}[A\mathbf{x} + b] = A\mathbb{E}[\mathbf{x}] + b \quad (2.96)$$

this is called the **linearity of expectation**.

For the covariance, we have

$$\text{Cov}[\mathbf{y}] = \text{Cov}[A\mathbf{x} + b] = A\Sigma A^T \quad (2.97)$$

## 2.9.2 General transformations

If  $X$  is a discrete rv, we can derive the pmf for  $y$  by simply summing up the probability mass for all the  $x$ s such that  $f(x) = y$ :

$$p_Y(y) = \sum_{x:g(x)=y} p_X(x) \quad (2.98)$$

If  $X$  is continuous, we cannot use Equation 2.98 since  $p_X(x)$  is a density, not a pmf, and we cannot sum up densities. Instead, we work with cdfs, and write

$$F_Y(y) = P(Y \leq y) = P(g(X) \leq y) = \int_{g(X) \leq y} f_X(x) dx \quad (2.99)$$

We can derive the pdf of  $Y$  by differentiating the cdf:

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right| \quad (2.100)$$

This is called **change of variables** formula. We leave the proof of this as an exercise.

For example, suppose  $X \sim U(1, 1)$ , and  $Y = X^2$ . Then  $p_Y(y) = \frac{1}{2}y^{-\frac{1}{2}}$ .

### 2.9.2.1 Multivariate change of variables \*

Let  $f$  be a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ , and let  $\mathbf{y} = f(\mathbf{x})$ . Then its Jacobian matrix  $\mathbf{J}$  is given by

$$\mathbf{J}_{\mathbf{x} \rightarrow \mathbf{y}} \triangleq \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \triangleq \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_n} \end{pmatrix} \quad (2.101)$$

$|\det(\mathbf{J})|$  measures how much a unit cube changes in volume when we apply  $f$ .

If  $f$  is an invertible mapping, we can define the pdf of the transformed variables using the Jacobian of the inverse mapping  $\mathbf{y} \rightarrow \mathbf{x}$ :

$$p_Y(\mathbf{y}) = p_X(\mathbf{x}) \left| \det \left( \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right) \right| = p_X(\mathbf{x}) |\det(\mathbf{J}_{\mathbf{y} \rightarrow \mathbf{x}})| \quad (2.102)$$

## 2.9.3 Central limit theorem

Given  $N$  random variables  $X_1, X_2, \dots, X_N$ , each variable is **independent and identically distributed**<sup>10</sup> (iid for short), and each has the same mean  $\mu$  and variance  $\sigma^2$ , then

$$\frac{\sum_{i=1}^n X_i - N\mu}{\sqrt{N}\sigma} \sim \mathcal{N}(0, 1) \quad (2.103)$$

this can also be written as

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{N}} \sim \mathcal{N}(0, 1) \quad , \text{ where } \bar{X} \triangleq \frac{1}{N} \sum_{i=1}^n X_i \quad (2.104)$$

<sup>10</sup> [http://en.wikipedia.org/wiki/Independent\\_identically\\_distributed](http://en.wikipedia.org/wiki/Independent_identically_distributed)

## 2.10 Monte Carlo approximation

In general, computing the distribution of a function of an rv using the change of variables formula can be difficult. One simple but powerful alternative is as follows. First we generate  $S$  samples from the distribution, call them  $x_1, \dots, x_S$ . (There are many ways to generate such samples; one popular method, for high dimensional distributions, is called Markov chain Monte Carlo or MCMC; this will be explained in Chapter TODO.) Given the samples, we can approximate the distribution of  $f(X)$  by using the empirical distribution of  $\{f(x_s)\}_{s=1}^S$ . This is called a **Monte Carlo approximation**<sup>11</sup>, named after a city in Europe known for its plush gambling casinos.

We can use Monte Carlo to approximate the expected value of any function of a random variable. We simply draw samples, and then compute the arithmetic mean of the function applied to the samples. This can be written as follows:

$$\mathbb{E}[g(X)] = \int g(x)p(x)dx \approx \frac{1}{S} \sum_{s=1}^S f(x_s) \quad (2.105)$$

where  $x_s \sim p(X)$ .

This is called **Monte Carlo integration**<sup>12</sup>, and has the advantage over numerical integration (which is based on evaluating the function at a fixed grid of points) that the function is only evaluated in places where there is non-negligible probability.

## 2.11 Information theory

### 2.11.1 Entropy

The entropy of a random variable  $X$  with distribution  $p$ , denoted by  $\mathbb{H}(X)$  or sometimes  $\mathbb{H}(p)$ , is a measure of its uncertainty. In particular, for a discrete variable with  $K$  states, it is defined by

$$\mathbb{H}(X) \triangleq - \sum_{k=1}^K p(X=k) \log_2 p(X=k) \quad (2.106)$$

Usually we use log base 2, in which case the units are called **bits**(short for binary digits). If we use log base  $e$ , the units are called **nats**.

The discrete distribution with maximum entropy is the uniform distribution (see Section XXX for a proof). Hence for a  $K$ -ary random variable, the entropy is maximized if  $p(x=k) = 1/K$ ; in this case,  $\mathbb{H}(X) = \log_2 K$ .

Conversely, the distribution with minimum entropy (which is zero) is any **delta-function** that puts all its mass on one state. Such a distribution has no uncertainty.

### 2.11.2 KL divergence

One way to measure the dissimilarity of two probability distributions,  $p$  and  $q$ , is known as the **Kullback-Leibler divergence**(**KL divergence**)or **relative entropy**. This is defined as follows:

$$\mathbb{KL}(P||Q) \triangleq \sum_x p(x) \log_2 \frac{p(x)}{q(x)} \quad (2.107)$$

<sup>11</sup> [http://en.wikipedia.org/wiki/Monte\\_Carlo\\_method](http://en.wikipedia.org/wiki/Monte_Carlo_method)

<sup>12</sup> [http://en.wikipedia.org/wiki/Monte\\_Carlo\\_integration](http://en.wikipedia.org/wiki/Monte_Carlo_integration)

where the sum gets replaced by an integral for pdfs<sup>13</sup>. The KL divergence is only defined if  $P$  and  $Q$  both sum to 1 and if  $q(x) = 0$  implies  $p(x) = 0$  for all  $x$  (absolute continuity). If the quantity  $0 \ln 0$  appears in the formula, it is interpreted as zero because  $\lim_{x \rightarrow 0} x \ln x$ . We can rewrite this as

$$\begin{aligned}\mathbb{KL}(p||q) &\triangleq \sum_x p(x) \log_2 p(x) - \sum_{k=1}^K p(x) \log_2 q(x) \\ &= \mathbb{H}(p) - \mathbb{H}(p, q)\end{aligned}\quad (2.108)$$

where  $\mathbb{H}(p, q)$  is called the **cross entropy**,

$$\mathbb{H}(p, q) = \sum_x p(x) \log_2 q(x) \quad (2.109)$$

One can show (Cover and Thomas 2006) that the cross entropy is the average number of bits needed to encode data coming from a source with distribution  $p$  when we use model  $q$  to define our codebook. Hence the regular entropy  $\mathbb{H}(p) = \mathbb{H}(p, p)$ , defined in section §2.11.1, is the expected number of bits if we use the true model, so the KL divergence is the difference between these. In other words, the KL divergence is the average number of *extra* bits needed to encode the data, due to the fact that we used distribution  $q$  to encode the data instead of the true distribution  $p$ .

The extra number of bits interpretation should make it clear that  $\mathbb{KL}(p||q) \geq 0$ , and that the KL is only equal to zero if  $q = p$ . We now give a proof of this important result.

**Theorem 2.1. (Information inequality)**  $\mathbb{KL}(p||q) \geq 0$  with equality iff  $p = q$ .

One important consequence of this result is that *the discrete distribution with the maximum entropy is the uniform distribution*.

### 2.11.3 Mutual information

**Definition 2.9. Mutual information** or **MI**, is defined as follows:

$$\begin{aligned}\mathbb{I}(X; Y) &\triangleq \mathbb{KL}(P(X, Y) || P(X)P(Y)) \\ &= \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)}\end{aligned}\quad (2.110)$$

We have  $\mathbb{I}(X; Y) \geq 0$  with equality if  $P(X, Y) = P(X)P(Y)$ . That is, the MI is zero if the variables are independent.

To gain insight into the meaning of MI, it helps to re-express it in terms of joint and conditional entropies. One can show that the above expression is equivalent to the following:

$$\mathbb{I}(X; Y) = \mathbb{H}(X) - \mathbb{H}(X|Y) \quad (2.111)$$

$$= \mathbb{H}(Y) - \mathbb{H}(Y|X) \quad (2.112)$$

$$= \mathbb{H}(X) + \mathbb{H}(Y) - \mathbb{H}(X, Y) \quad (2.113)$$

$$= \mathbb{H}(X, Y) - \mathbb{H}(X|Y) - \mathbb{H}(Y|X) \quad (2.114)$$

where  $\mathbb{H}(X)$  and  $\mathbb{H}(Y)$  are the **marginal entropies**,  $\mathbb{H}(X|Y)$  and  $\mathbb{H}(Y|X)$  are the **conditional entropies**, and  $\mathbb{H}(X, Y)$  is the **joint entropy** of  $X$  and  $Y$ , see Fig. 2.10<sup>14</sup>.

Intuitively, we can interpret the MI between  $X$  and  $Y$  as the reduction in uncertainty about  $X$  after observing  $Y$ , or, by symmetry, the reduction in uncertainty about  $Y$  after observing  $X$ .

<sup>13</sup> The KL divergence is not a distance, since it is asymmetric. One symmetric version of the KL divergence is the **Jensen-Shannon divergence**, defined as  $JS(p_1, p_2) = 0.5\mathbb{KL}(p_1||q) + 0.5\mathbb{KL}(p_2||q)$ , where  $q = 0.5p_1 + 0.5p_2$

<sup>14</sup> [http://en.wikipedia.org/wiki/Mutual\\_information](http://en.wikipedia.org/wiki/Mutual_information)

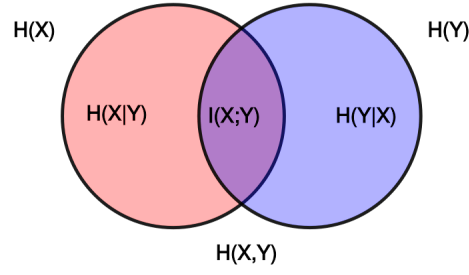


Fig. 2.10: Individual  $\mathbb{H}(X)$ ,  $\mathbb{H}(Y)$ , joint  $\mathbb{H}(X, Y)$ , and conditional entropies for a pair of correlated subsystems  $X, Y$  with mutual information  $\mathbb{I}(X; Y)$ .

A quantity which is closely related to MI is the **pointwise mutual information** or **PMI**. For two events (not random variables)  $x$  and  $y$ , this is defined as

$$PMI(x, y) \triangleq \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)} \quad (2.115)$$

This measures the discrepancy between these events occurring together compared to what would be expected by chance. Clearly the MI of  $X$  and  $Y$  is just the expected value of the PMI. Interestingly, we can rewrite the PMI as follows:

$$PMI(x, y) = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)} \quad (2.116)$$

This is the amount we learn from updating the prior  $p(x)$  into the posterior  $p(x|y)$ , or equivalently, updating the prior  $p(y)$  into the posterior  $p(y|x)$ .



## Chapter 3

### Generative models for discrete data

#### 3.1 Generative classifier

$$p(y = c | \mathbf{x}, \boldsymbol{\theta}) = \frac{p(y = c | \boldsymbol{\theta}) p(\mathbf{x} | y = c, \boldsymbol{\theta})}{\sum_{c'} p(y = c' | \boldsymbol{\theta}) p(\mathbf{x} | y = c', \boldsymbol{\theta})} \quad (3.1)$$

This is called a **generative classifier**, since it specifies how to generate the data using the **class conditional density**  $p(\mathbf{x} | y = c)$  and the class prior  $p(y = c)$ . An alternative approach is to directly fit the class posterior,  $p(y = c | \mathbf{x})$ ; this is known as a **discriminative classifier**.

#### 3.2 Bayesian concept learning

Psychological research has shown that people can learn concepts from positive examples alone (Xu and Tenenbaum 2007).

We can think of learning the meaning of a word as equivalent to **concept learning**, which in turn is equivalent to binary classification. To see this, define  $f(\mathbf{x}) = 1$  if  $\mathbf{x}$  is an example of the concept  $C$ , and  $f(\mathbf{x}) = 0$  otherwise. Then the goal is to learn the indicator function  $f$ , which just defines which elements are in the set  $C$ .

##### 3.2.1 Likelihood

$$p(\mathcal{D} | h) \triangleq \left( \frac{1}{\text{size}(h)} \right)^N = \left( \frac{1}{|h|} \right)^N \quad (3.2)$$

This crucial equation embodies what Tenenbaum calls the **size principle**, which means the model favours the simplest (smallest) hypothesis consistent with the data. This is more commonly known as **Occams razor**<sup>15</sup>.

##### 3.2.2 Prior

The prior is decided by human, not machines, so it is subjective. The subjectivity of the prior is controversial. For example, that a child and a math professor will reach different answers. In fact, they presumably not only have different priors, but also different hypothesis spaces. However, we can finesse that by defining the hypothesis space of the child and the math professor to be the same, and then setting the child's prior weight to be zero on certain advanced concepts. Thus there is no sharp distinction between the prior and the hypothesis space.

However, the prior is the mechanism by which background knowledge can be brought to bear on a problem. Without this, rapid learning (i.e., from small samples sizes) is impossible.

##### 3.2.3 Posterior

The posterior is simply the likelihood times the prior, normalized.

---

<sup>15</sup> [http://en.wikipedia.org/wiki/Occam%27s\\_razor](http://en.wikipedia.org/wiki/Occam%27s_razor)

$$p(h|\mathcal{D}) \triangleq \frac{p(\mathcal{D}|h)p(h)}{\sum_{h' \in \mathcal{H}} p(\mathcal{D}|h')p(h')} = \frac{\mathbb{I}(\mathcal{D} \in h)p(h)}{\sum_{h' \in \mathcal{H}} \mathbb{I}(\mathcal{D} \in h')p(h')} \quad (3.3)$$

where  $\mathbb{I}(\mathcal{D} \in h)p(h)$  is 1 **iff**(iff and only if) all the data are in the extension of the hypothesis  $h$ .

In general, when we have enough data, the posterior  $p(h|\mathcal{D})$  becomes peaked on a single concept, namely the MAP estimate, i.e.,

$$p(h|\mathcal{D}) \rightarrow \hat{h}^{MAP} \quad (3.4)$$

where  $\hat{h}^{MAP}$  is the posterior mode,

$$\begin{aligned} \hat{h}^{MAP} &\triangleq \arg \max_h p(h|\mathcal{D}) = \arg \max_h p(\mathcal{D}|h)p(h) \\ &= \arg \max_h [\log p(\mathcal{D}|h) + \log p(h)] \end{aligned} \quad (3.5)$$

Since the likelihood term depends exponentially on  $N$ , and the prior stays constant, as we get more and more data, the MAP estimate converges towards the **maximum likelihood estimate** or **MLE**:

$$\hat{h}^{MLE} \triangleq \arg \max_h p(\mathcal{D}|h) = \arg \max_h \log p(\mathcal{D}|h) \quad (3.6)$$

In other words, if we have enough data, we see that the **data overwhelms the prior**.

### 3.2.4 Posterior predictive distribution

The concept of **posterior predictive distribution**<sup>16</sup> is normally used in a Bayesian context, where it makes use of the entire posterior distribution of the parameters given the observed data to yield a probability distribution over an interval rather than simply a point estimate.

$$p(\tilde{x}|\mathcal{D}) \triangleq \mathbb{E}_{h|\mathcal{D}}[p(\tilde{x}|h)] = \begin{cases} \sum_h p(\tilde{x}|h)p(h|\mathcal{D}) \\ \int p(\tilde{x}|h)p(h|\mathcal{D})dh \end{cases} \quad (3.7)$$

This is just a weighted average of the predictions of each individual hypothesis and is called **Bayes model averaging**(Hoeting et al. 1999).

## 3.3 The beta-binomial model

### 3.3.1 Likelihood

Given  $X \sim \text{Bin}(\theta)$ , the likelihood of  $\mathcal{D}$  is given by

$$p(\mathcal{D}|\theta) = \text{Bin}(N_1|N, \theta) \quad (3.8)$$

### 3.3.2 Prior

$$\text{Beta}(\theta|a, b) \propto \theta^{a-1}(1-\theta)^{b-1} \quad (3.9)$$

The parameters of the prior are called **hyper-parameters**.

<sup>16</sup> [http://en.wikipedia.org/wiki/Posterior\\_predictive\\_distribution](http://en.wikipedia.org/wiki/Posterior_predictive_distribution)

### 3.3.3 Posterior

$$\begin{aligned} p(\theta|\mathcal{D}) &\propto \text{Bin}(N_1|N_1 + N_0, \theta) \text{Beta}(\theta|a, b) \\ &= \text{Beta}(\theta|N_1 + a, N_0 + b) \end{aligned} \quad (3.10)$$

Note that updating the posterior sequentially is equivalent to updating in a single batch. To see this, suppose we have two data sets  $\mathcal{D}_a$  and  $\mathcal{D}_b$  with sufficient statistics  $N_1^a, N_0^a$  and  $N_1^b, N_0^b$ . Let  $N_1 = N_1^a + N_1^b$  and  $N_0 = N_0^a + N_0^b$  be the sufficient statistics of the combined datasets. In batch mode we have

$$\begin{aligned} p(\theta|\mathcal{D}_a, \mathcal{D}_b) &= p(\theta, \mathcal{D}_b|\mathcal{D}_a) p(\mathcal{D}_b) \\ &\propto p(\theta, \mathcal{D}_b|\mathcal{D}_a) \\ &= p(\mathcal{D}_b, \theta|\mathcal{D}_a) \\ &= p(\mathcal{D}_b|\theta) p(\theta|\mathcal{D}_a) \\ &\text{Combine Equation 3.10 and 2.44} \\ &= \text{Bin}(N_1^b|\theta, N_1^b + N_0^b) \text{Beta}(\theta|N_1^a + a, N_0^a + b) \\ &= \text{Beta}(\theta|N_1^a + N_1^b + a, N_0^a + N_0^b + b) \end{aligned}$$

This makes Bayesian inference particularly well-suited to **online learning**, as we will see later.

#### 3.3.3.1 Posterior mean and mode

From Table 2.7, the posterior mean is given by

$$\bar{\theta} = \frac{a + N_1}{a + b + N} \quad (3.11)$$

The mode is given by

$$\hat{\theta}_{MAP} = \frac{a + N_1 - 1}{a + b + N - 2} \quad (3.12)$$

If we use a uniform prior, then the MAP estimate reduces to the MLE,

$$\hat{\theta}_{MLE} = \frac{N_1}{N} \quad (3.13)$$

We will now show that the posterior mean is convex combination of the prior mean and the MLE, which captures the notion that the posterior is a compromise between what we previously believed and what the data is telling us.

#### 3.3.3.2 Posterior variance

The mean and mode are point estimates, but it is useful to know how much we can trust them. The variance of the posterior is one way to measure this. The variance of the Beta posterior is given by

$$\text{var}(\theta|\mathcal{D}) = \frac{(a + N_1)(b + N_0)}{(a + N_1 + b + N_0)^2(a + N_1 + b + N_0 + 1)} \quad (3.14)$$

We can simplify this formidable expression in the case that  $N \gg a, b$ , to get

$$\text{var}(\theta|\mathcal{D}) \approx \frac{N_1 N_0}{N N N} = \frac{\hat{\theta}_{MLE}(1 - \hat{\theta}_{MLE})}{N} \quad (3.15)$$

### 3.3.4 Posterior predictive distribution

So far, we have been focusing on inference of the unknown parameter(s). Let us now turn our attention to prediction of future observable data.

Consider predicting the probability of heads in a single future trial under a  $\text{Beta}(a, b)$  posterior. We have

$$\begin{aligned} p(\tilde{x}|\mathcal{D}) &= \int_0^1 p(\tilde{x}|\theta) p(\theta|\mathcal{D}) d\theta \\ &= \int_0^1 \theta \text{Beta}(\theta|a, b) d\theta \\ &= \mathbb{E}[\theta|\mathcal{D}] = \frac{a}{a+b} \end{aligned} \quad (3.16)$$

#### 3.3.4.1 Overfitting and the black swan paradox

Let us now derive a simple Bayesian solution to the problem. We will use a uniform prior, so  $a = b = 1$ . In this case, plugging in the posterior mean gives **Laplace's rule of succession**

$$p(\tilde{x}|\mathcal{D}) = \frac{N_1 + 1}{N_0 + N_1 + 1} \quad (3.17)$$

This justifies the common practice of adding 1 to the empirical counts, normalizing and then plugging them in, a technique known as **add-one smoothing**. (Note that plugging in the MAP parameters would not have this smoothing effect, since the mode becomes the MLE if  $a = b = 1$ , see Section 3.3.3.1.)

#### 3.3.4.2 Predicting the outcome of multiple future trials

Suppose now we were interested in predicting the number of heads,  $\tilde{x}$ , in  $M$  future trials. This is given by

$$p(\tilde{x}|\mathcal{D}) = \int_0^1 \text{Bin}(\tilde{x}|M, \theta) \text{Beta}(\theta|a, b) d\theta \quad (3.18)$$

$$= \binom{M}{\tilde{x}} \frac{1}{B(a, b)} \int_0^1 \theta^{\tilde{x}} (1 - \theta)^{M - \tilde{x}} \theta^{a-1} (1 - \theta)^{b-1} d\theta \quad (3.19)$$

We recognize the integral as the normalization constant for a  $\text{Beta}(a + \tilde{x}, M\tilde{x} + b)$  distribution. Hence

$$\int_0^1 \theta^{\tilde{x}} (1 - \theta)^{M - \tilde{x}} \theta^{a-1} (1 - \theta)^{b-1} d\theta = B(\tilde{x} + a, M - \tilde{x} + b) \quad (3.20)$$

Thus we find that the posterior predictive is given by the following, known as the (compound) **beta-binomial distribution**:

$$Bb(x|a, b, M) \triangleq \binom{M}{x} \frac{B(x + a, M - x + b)}{B(a, b)} \quad (3.21)$$

This distribution has the following mean and variance

$$\text{mean} = M \frac{a}{a+b}, \text{var} = \frac{Mab}{(a+b)^2} \frac{a+b+M}{a+b+1} \quad (3.22)$$

This process is illustrated in Figure 3.1. We start with a  $\text{Beta}(2, 2)$  prior, and plot the posterior predictive density after seeing  $N_1 = 3$  heads and  $N_0 = 17$  tails. Figure 3.1(b) plots a plug-in approximation using a MAP estimate. We see that the Bayesian prediction has longer tails, spreading its probability mass more widely, and is therefore less prone to overfitting and blackswan type paradoxes.

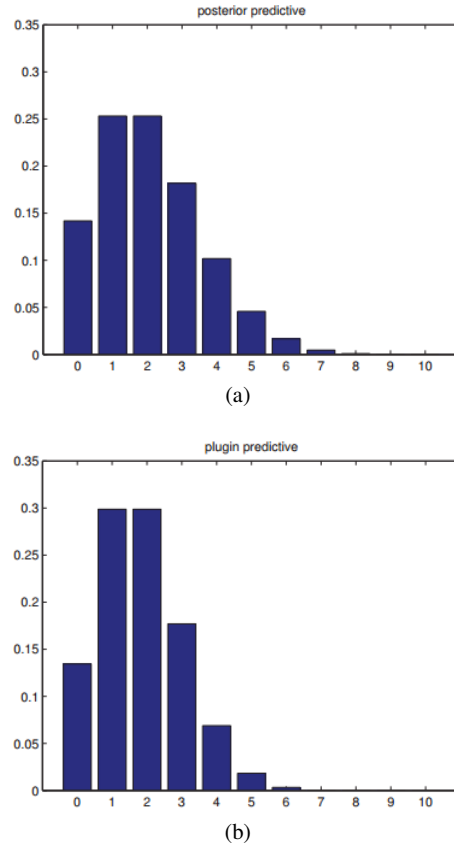


Fig. 3.1: (a) Posterior predictive distributions after seeing  $N_1 = 3, N_0 = 17$ . (b) MAP estimation.

### 3.4 The Dirichlet-multinomial model

In the previous section, we discussed how to infer the probability that a coin comes up heads. In this section, we generalize these results to infer the probability that a dice with  $K$  sides comes up as face  $k$ .

#### 3.4.1 Likelihood

Suppose we observe  $N$  dice rolls,  $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ , where  $x_i \in \{1, 2, \dots, K\}$ . The likelihood has the form

$$p(\mathcal{D}|\theta) = \binom{N}{N_1 \dots N_K} \prod_{k=1}^K \theta_k^{N_k} \quad \text{where } N_k = \sum_{i=1}^N \mathbb{I}(y_i = k) \quad (3.23)$$

almost the same as Equation 2.46.

#### 3.4.2 Prior

$$\text{Dir}(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \mathbb{I}(\theta \in S_K) \quad (3.24)$$

### 3.4.3 Posterior

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (3.25)$$

$$\propto \prod_{k=1}^K \theta_k^{N_k} \theta_k^{\alpha_k-1} = \prod_{k=1}^K \theta_k^{N_k+\alpha_k-1} \quad (3.26)$$

$$= \text{Dir}(\boldsymbol{\theta}|\alpha_1 + N_1, \dots, \alpha_K + N_K) \quad (3.27)$$

From Equation 2.94, the MAP estimate is given by

$$\hat{\theta}_k = \frac{N_k + \alpha_k - 1}{N + \alpha_0 - K} \quad (3.28)$$

If we use a uniform prior,  $\alpha_k = 1$ , we recover the MLE:

$$\hat{\theta}_k = \frac{N_k}{N} \quad (3.29)$$

### 3.4.4 Posterior predictive distribution

The posterior predictive distribution for a single multinoulli trial is given by the following expression:

$$p(X = j|\mathcal{D}) = \int p(X = j|\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta} \quad (3.30)$$

$$= \int p(X = j|\theta_j) \left[ \int p(\boldsymbol{\theta}_{-j}, \theta_j|\mathcal{D})d\boldsymbol{\theta}_{-j} \right] d\theta_j \quad (3.31)$$

$$= \int \theta_j p(\theta_j|\mathcal{D})d\theta_j = \mathbb{E}[\theta_j|\mathcal{D}] = \frac{\alpha_j + N_j}{\alpha_0 + N} \quad (3.32)$$

where  $\boldsymbol{\theta}_{-j}$  are all the components of  $\boldsymbol{\theta}$  except  $\theta_j$ .

The above expression avoids the zero-count problem. In fact, this form of Bayesian smoothing is even more important in the multinomial case than the binary case, since the likelihood of data sparsity increases once we start partitioning the data into many categories.

## 3.5 Naive Bayes classifiers

Assume the features are **conditionally independent** given the class label, then the class conditional density has the following form

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \prod_{j=1}^D p(x_j|y = c, \boldsymbol{\theta}_{jc}) \quad (3.33)$$

The resulting model is called a **naive Bayes classifier**(NBC).

The form of the class-conditional density depends on the type of each feature. We give some possibilities below:

- In the case of real-valued features, we can use the Gaussian distribution:  $p(\mathbf{x}|y, \boldsymbol{\theta}) = \prod_{j=1}^D \mathcal{N}(x_j|\mu_{jc}, \sigma_{jc}^2)$ , where  $\mu_{jc}$  is the mean of feature  $j$  in objects of class  $c$ , and  $\sigma_{jc}^2$  is its variance.
- In the case of binary features,  $x_j \in \{0, 1\}$ , we can use the Bernoulli distribution:  $p(\mathbf{x}|y, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Ber}(x_j|\mu_{jc})$ , where  $\mu_{jc}$  is the probability that feature  $j$  occurs in class  $c$ . This is sometimes called the **multivariate Bernoulli naive Bayes** model. We will see an application of this below.
- In the case of categorical features,  $x_j \in \{a_{j1}, a_{j2}, \dots, a_{jS_j}\}$ , we can use the multinoulli distribution:  $p(\mathbf{x}|y, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Cat}(x_j|\boldsymbol{\mu}_{jc})$ , where  $\boldsymbol{\mu}_{jc}$  is a histogram over the  $K$  possible values for  $x_j$  in class  $c$ .

Obviously we can handle other kinds of features, or use different distributional assumptions. Also, it is easy to mix and match features of different types.

### 3.5.1 Optimization

We now discuss how to train a naive Bayes classifier. This usually means computing the MLE or the MAP estimate for the parameters. However, we will also discuss how to compute the full posterior,  $p(\boldsymbol{\theta}|\mathcal{D})$ .

#### 3.5.1.1 MLE for NBC

The probability for a single data case is given by

$$\begin{aligned} p(\mathbf{x}_i, y_i | \boldsymbol{\theta}) &= p(y_i | \boldsymbol{\pi}) \prod_j p(x_{ij} | \boldsymbol{\theta}_j) \\ &= \prod_c \pi_c^{\mathbb{I}(y_i=c)} \prod_j \prod_c p(x_{ij} | \boldsymbol{\theta}_{jc})^{\mathbb{I}(y_i=c)} \end{aligned} \quad (3.34)$$

Hence the log-likelihood is given by

$$p(\mathcal{D} | \boldsymbol{\theta}) = \sum_{c=1}^C N_c \log \pi_c + \sum_{j=1}^D \sum_{c=1}^C \sum_{i: y_i=c} \log p(x_{ij} | \boldsymbol{\theta}_{jc}) \quad (3.35)$$

where  $N_c \triangleq \sum_i \mathbb{I}(y_i = c)$  is the number of feature vectors in class  $c$ .

We see that this expression decomposes into a series of terms, one concerning  $\boldsymbol{\pi}$ , and  $DC$  terms containing the  $\boldsymbol{\theta}_{jc}$ s. Hence we can optimize all these parameters separately.

From Equation 3.29, the MLE for the class prior is given by

$$\hat{\pi}_c = \frac{N_c}{N} \quad (3.36)$$

The MLE for  $\boldsymbol{\theta}_{jc}$ s depends on the type of distribution we choose to use for each feature.

In the case of binary features,  $x_j \in \{0, 1\}$ ,  $x_j | y = c \sim \text{Ber}(\boldsymbol{\theta}_{jc})$ , hence

$$\hat{\boldsymbol{\theta}}_{jc} = \frac{N_{jc}}{N_c} \quad (3.37)$$

where  $N_{jc} \triangleq \sum_{i: y_i=c} \mathbb{I}(x_{ij} = 1)$  is the number that feature  $j$  occurs in class  $c$ .

In the case of categorical features,  $x_j \in \{a_{j1}, a_{j2}, \dots, a_{jS_j}\}$ ,  $x_j | y = c \sim \text{Cat}(\boldsymbol{\theta}_{jc})$ , hence

$$\hat{\boldsymbol{\theta}}_{jc} = \left( \frac{N_{j1c}}{N_c}, \frac{N_{j2c}}{N_c}, \dots, \frac{N_{jS_j c}}{N_c} \right)^T \quad (3.38)$$

where  $N_{jkc} \triangleq \sum_{i=1}^N \mathbb{I}(x_{ij} = a_{jk}, y_i = c)$  is the number that feature  $x_j = a_{jk}$  occurs in class  $c$ .

#### 3.5.1.2 Bayesian naive Bayes

Use a  $\text{Dir}(\boldsymbol{\alpha})$  prior for  $\boldsymbol{\pi}$ .

In the case of binary features, use a  $\text{Beta}(\beta_0, \beta_1)$  prior for each  $\boldsymbol{\theta}_{jc}$ ; in the case of categorical features, use a  $\text{Dir}(\boldsymbol{\alpha})$  prior for each  $\boldsymbol{\theta}_{jc}$ . Often we just take  $\boldsymbol{\alpha} = \mathbf{1}$  and  $\boldsymbol{\beta} = \mathbf{1}$ , corresponding to **add-one** or **Laplace smoothing**.

### 3.5.2 Using the model for prediction

The goal is to compute

$$\begin{aligned} y = f(\mathbf{x}) &= \arg \max_c P(y = c | \mathbf{x}, \boldsymbol{\theta}) \\ &= P(y = c | \boldsymbol{\theta}) \prod_{j=1}^D P(x_j | y = c, \boldsymbol{\theta}) \end{aligned} \quad (3.39)$$

We can estimate parameters using MLE or MAP, then the posterior predictive density is obtained by simply plugging in the parameters  $\bar{\boldsymbol{\theta}}$ (MLE) or  $\hat{\boldsymbol{\theta}}$ (MAP).

Or we can use BMA, just integrate out the unknown parameters.

### 3.5.3 The log-sum-exp trick

when using generative classifiers of any kind, computing the posterior over class labels using Equation 3.1 can fail due to **numerical underflow**. The problem is that  $p(\mathbf{x} | y = c)$  is often a very small number, especially if  $\mathbf{x}$  is a high-dimensional vector. This is because we require that  $\sum_{\mathbf{x}} p(\mathbf{x} | y) = 1$ , so the probability of observing any particular high-dimensional vector is small. The obvious solution is to take logs when applying Bayes rule, as follows:

$$\log p(y = c | \mathbf{x}, \boldsymbol{\theta}) = b_c - \log \left( \sum_{c'} e^{b_{c'}} \right) \quad (3.40)$$

where  $b_c \triangleq \log p(\mathbf{x} | y = c, \boldsymbol{\theta}) + \log p(y = c | \boldsymbol{\theta})$ .

We can factor out the largest term, and just represent the remaining numbers relative to that. For example,

$$\begin{aligned} \log(e^{-120} + e^{-121}) &= \log(e^{-120}(1 + e^{-1})) \\ &= \log(1 + e^{-1}) - 120 \end{aligned} \quad (3.41)$$

In general, we have

$$\sum_c e^{b_c} = \log \left[ \left( \sum_c e^{b_c - B} \right) e^B \right] = \log \left( \sum_c e^{b_c - B} \right) + B \quad (3.42)$$

where  $B \triangleq \max\{b_c\}$ .

This is called the **log-sum-exp** trick, and is widely used.

### 3.5.4 Feature selection using mutual information

Since an NBC is fitting a joint distribution over potentially many features, it can suffer from overfitting. In addition, the run-time cost is  $O(D)$ , which may be too high for some applications.

One common approach to tackling both of these problems is to perform **feature selection**, to remove irrelevant features that do not help much with the classification problem. The simplest approach to feature selection is to evaluate the relevance of each feature separately, and then take the top  $K$ , where  $K$  is chosen based on some tradeoff between accuracy and complexity. This approach is known as **variable ranking**, **filtering**, or **screening**.

One way to measure relevance is to use mutual information (Section 2.11.3) between feature  $X_j$  and the class label  $Y$

$$\mathbb{I}(X_j, Y) = \sum_{x_j} \sum_y p(x_j, y) \log \frac{p(x_j, y)}{p(x_j)p(y)} \quad (3.43)$$

If the features are binary, it is easy to show that the MI can be computed as follows



$$\mathbb{I}_j = \sum_c \left[ \theta_{jc} \pi_c \log \frac{\theta_{jc}}{\theta_j} + (1 - \theta_{jc}) \pi_c \log \frac{1 - \theta_{jc}}{1 - \theta_j} \right] \quad (3.44)$$

where  $\pi_c = p(y = c)$ ,  $\theta_{jc} = p(x_j = 1 | y = c)$ , and  $\theta_j = p(x_j = 1) = \sum_c \pi_c \theta_{jc}$ .

### 3.5.5 Classifying documents using bag of words

**Document classification** is the problem of classifying text documents into different categories.

#### 3.5.5.1 Bernoulli product model

One simple approach is to represent each document as a binary vector, which records whether each word is present or not, so  $x_{ij} = 1$  iff word  $j$  occurs in document  $i$ , otherwise  $x_{ij} = 0$ . We can then use the following class conditional density:

$$\begin{aligned} p(\mathbf{x}_i | y_i = c, \boldsymbol{\theta}) &= \prod_{j=1}^D \text{Ber}(x_{ij} | \theta_{jc}) \\ &= \prod_{j=1}^D \theta_{jc}^{x_{ij}} (1 - \theta_{jc})^{1-x_{ij}} \end{aligned} \quad (3.45)$$

This is called the **Bernoulli product model**, or the **binary independence model**.

#### 3.5.5.2 Multinomial document classifier

However, ignoring the number of times each word occurs in a document loses some information (McCallum and Nigam 1998). A more accurate representation counts the number of occurrences of each word. Specifically, let  $\mathbf{x}_i$  be a vector of counts for document  $i$ , so  $x_{ij} \in \{0, 1, \dots, N_i\}$ , where  $N_i$  is the number of terms in document  $i$  (so  $\sum_{j=1}^D x_{ij} = N_i$ ).

For the class conditional densities, we can use a multinomial distribution:

$$p(\mathbf{x}_i | y_i = c, \boldsymbol{\theta}) = \text{Mu}(\mathbf{x}_i | N_i, \boldsymbol{\theta}_c) = \frac{N_i!}{\prod_{j=1}^D x_{ij}!} \prod_{j=1}^D \theta_{jc}^{x_{ij}} \quad (3.46)$$

where we have implicitly assumed that the document length  $N_i$  is independent of the class. Here  $\theta_{jc}$  is the probability of generating word  $j$  in documents of class  $c$ ; these parameters satisfy the constraint that  $\sum_{j=1}^D \theta_{jc} = 1$  for each class  $c$ .

Although the multinomial classifier is easy to train and easy to use at test time, it does not work particularly well for document classification. One reason for this is that it does not take into account the **burstiness** of word usage. This refers to the phenomenon that most words never appear in any given document, but if they do appear once, they are likely to appear more than once, i.e., words occur in bursts.

The multinomial model cannot capture the burstiness phenomenon. To see why, note that Equation 3.46 has the form  $\theta_{jc}^{x_{ij}}$ , and since  $\theta_{jc} \ll 1$  for rare words, it becomes increasingly unlikely to generate many of them. For more frequent words, the decay rate is not as fast. To see why intuitively, note that the most frequent words are function words which are not specific to the class, such as and, the, and but; the chance of the word and occurring is pretty much the same no matter how many time it has previously occurred (modulo document length), so the independence assumption is more reasonable for common words. However, since rare words are the ones that matter most for classification purposes, these are the ones we want to model the most carefully.

### 3.5.5.3 DCM model

Various ad hoc heuristics have been proposed to improve the performance of the multinomial document classifier (Rennie et al. 2003). We now present an alternative class conditional density that performs as well as these ad hoc methods, yet is probabilistically sound (Madsen et al. 2005).

Suppose we simply replace the multinomial class conditional density with the **Dirichlet Compound Multinomial** or **DCM** density, defined as follows:

$$\begin{aligned} p(\mathbf{x}_i|y_i = c, \boldsymbol{\alpha}) &= \int \text{Mu}(\mathbf{x}_i|N_i, \boldsymbol{\theta}_c) \text{Dir}(\boldsymbol{\theta}_c|\boldsymbol{\alpha}_c) \\ &= \frac{N_i!}{\prod_{j=1}^D x_{ij}!} \prod_{j=1}^D \frac{B(\mathbf{x}_i + \boldsymbol{\alpha}_c)}{B(\boldsymbol{\alpha}_c)} \end{aligned} \quad (3.47)$$

(This equation is derived in Equation TODO.) Surprisingly this simple change is all that is needed to capture the burstiness phenomenon. The intuitive reason for this is as follows: After seeing one occurrence of a word, say word<sub>j</sub>, the posterior counts on j gets updated, making another occurrence of word<sub>j</sub> more likely. By contrast, if j is fixed, then the occurrences of each word are independent. The multinomial model corresponds to drawing a ball from an urn with K colors of ball, recording its color, and then replacing it. By contrast, the DCM model corresponds to drawing a ball, recording its color, and then replacing it with one additional copy; this is called the **Polya urn**.

Using the DCM as the class conditional density gives much better results than using the multinomial, and has performance comparable to state of the art methods, as described in (Madsen et al. 2005). The only disadvantage is that fitting the DCM model is more complex; see (Minka 2000e; Elkan 2006) for the details.

## Chapter 4

### Gaussian Models

The Gaussian, also known as the **normal distribution**, is a widely used model for the **distribution of continuous variables**. In the case of a single variable  $x$ , the Gaussian distribution can be written in the form

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} \quad (4.1)$$

where the  $\mu$  is the mean and  $\sigma^2$  is the variance. Now we will discuss the **multivariate Gaussian** or **multivariate normal (MVN)**, which is the most widely used joint probability density function for continuous variables. It will form the basis for many of the models.

#### 4.1 Basics

Recall from Section 2.8.2, for a  $D$ -dimensional vector  $x$ , the multivariate Gaussian distribution (MVN), i.e. the **pdf** takes the form:

$$\mathcal{N}(x|\mu, \Sigma) \triangleq \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left[-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right] \quad (4.2)$$

where  $\mu$  is a  $D$ -dimensional mean vector,  $\Sigma$  is a  $D \times D$  covariance matrix, and  $|\Sigma|$  denotes the determinant of  $\Sigma$ .

The Gaussian distribution arises in many different contexts and can be motivated from variety of different perspectives, such as the distribution that maximizes the entropy. Another situation in which the Gaussian distribution arises is when we consider the sum of multiple random variables. The *central limit theorem* (due to Laplace), tells us that, subject to certain mild conditions, the sum of a set of random variables, which is of course itself a random variable, has a distribution that becomes increasingly Gaussian as the number of terms in the sum increases (Walker, 1969).

The expression inside the exponent, which is the functional dependence of the Gaussian on  $x$ , is through the quadratic form

$$\Delta^2 = (x-\mu)^T \Sigma^{-1} (x-\mu) \quad (4.3)$$

The quantity  $\Delta$  is the *Mahalanobis distance* between a data vector  $x$  and the mean vector  $\mu$  and reduces to the Euclidean distance when  $\Sigma$  is the identity matrix. We can gain a better understanding of this quantity by performing an **eigendecomposition** of  $\Sigma$ . That is, we write  $\Sigma = U \Lambda U^T$ , where  $U$  is an orthonormal matrix of eigenvectors satisfying  $U^T U = I$ , and  $\Lambda$  is a diagonal matrix of eigenvalues. Using the eigendecomposition, we have that

$$\Sigma^{-1} = U^{-T} \Lambda^{-1} U^{-1} = U \Lambda^{-1} U^T = \sum_{i=1}^D \frac{1}{\lambda_i} u_i u_i^T \quad (4.4)$$

where  $u_i$  is the  $i$ 'th column of  $U$ , containing the  $i$ 'th eigenvector. Hence we can rewrite the Mahalanobis distance as follows:

$$(x-\mu)^T \Sigma^{-1} (x-\mu) = (x-\mu)^T \left( \sum_{i=1}^D \frac{1}{\lambda_i} u_i u_i^T \right) (x-\mu) \quad (4.5)$$

$$= \sum_{i=1}^D \frac{1}{\lambda_i} (x-\mu)^T u_i u_i^T (x-\mu) \quad (4.6)$$

$$= \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \quad (4.7)$$

where  $y_i \triangleq u_i^T (x-\mu)$ . Recall that the equation for an ellipse in 2d is

$$\frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2} = 1 \quad (4.8)$$

Hence we see that the contours of equal probability density of a Gaussian lie along ellipses. This is illustrated in Figure 4.1. The eigenvectors determine the orientation of the ellipse, and the eigenvalues determine how elongated it is.

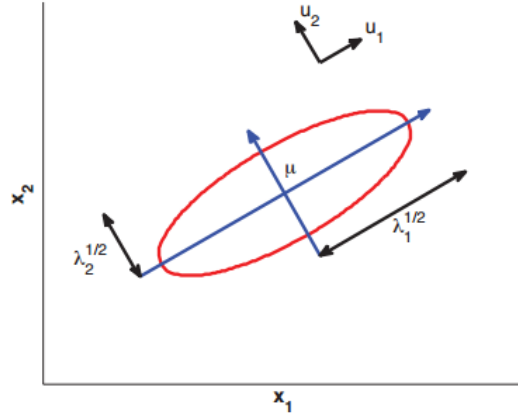


Fig. 4.1: Visualization of a 2 dimensional Gaussian density. The major and minor axes of the ellipse are defined by the first two eigenvectors of the covariance matrix, namely  $u_1$  and  $u_2$ . Based on Figure 2.7 of (Bishop 2006a)

In general, we see that the Mahalanobis distance corresponds to Euclidean distance in a transformed coordinate system, where we shift by  $\mu$  and rotate by  $U$ .

#### 4.1.1 proof

*Proof.* First of all, we note that the matrix  $\Sigma$  can be taken to be symmetric, without loss of generality. Now we consider the eigenvector equation for the covariance matrix

$$\Sigma \mu_i = \lambda_i \mu_i \quad (4.9)$$

where  $i = 1, \dots, D$ . Because  $\Sigma$  is a real, symmetric matrix its eigenvalues will be real, and its eigenvectors can be chosen to form an orthonormal set, so that

$$\mu_i^T \mu_j = \mathcal{I}_{ij} \quad (4.10)$$

where  $\mathcal{I}_{ij}$  is the  $i, j$  element of the identity matrix and satisfies

$$\mathcal{I}_{ij} = \begin{cases} 1, & \text{if } i=j \\ 0, & \text{otherwise.} \end{cases} \quad (4.11)$$

The covariance matrix  $\Sigma$  can be expressed as an equation in terms of its eigenvectors in the form

$$\Sigma = U \Lambda U^T \quad (4.12)$$

$$= \sum_{i=1}^D \lambda_i \mu_i \mu_i^T \quad (4.13)$$

where  $u_i$  is the  $i$ 'th column of  $U$ , containing the  $i$ 'th eigenvector. Similarly, the inverse covariance matrix  $\Sigma^{-1}$  can be expressed as

$$\Sigma^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad (4.14)$$

Substituting 4.14 into 4.3 ,the quadratic form becomes

$$\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \quad (4.15)$$

where we have defined

$$y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu}) \quad (4.16)$$

We can interpret  $\{y_i\}$  as a new coordinate system defined by the orthonormal vectors  $\mathbf{u}_i$  that are shifted and rotated with respect to the original  $x_i$  coordinates. Forming the vector  $\mathbf{y} = (y_1, \dots, y_D)^T$ , we have

$$\mathbf{y} = \mathbf{U}(\mathbf{x} - \boldsymbol{\mu}) \quad (4.17)$$

where

$$\mathbf{U} = [\mathbf{u}_1^T \ \mathbf{u}_2^T \ \dots \ \mathbf{u}_n^T] \quad (4.18)$$

$\mathbf{U}$  is an *orthogonal* matrix. A matrix whose eigenvalue are strictly positive is said to be *positive definite*, and if all the eigenvalues are nonnegative, then the covariance matrix is said to be *positive semidefinite*.

Now consider the form of Gaussian distribution in the new coordinate system defined by the  $y_i$ . In going from  $\mathbf{x}$  to  $\mathbf{y}$  coordinate system, we have a *Jacobian matrix*  $\mathbf{J}$  with elements given by

$$J_{ij} = \frac{\partial x_i}{\partial y_j} = U_{ji} \quad (4.19)$$

$$\mathbf{J} = \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \quad (4.20)$$

where  $U_{ji}$  are the elements of the matrix  $\mathbf{U}^T$ . Using the orthonormality property of the matrix  $\mathbf{U}$ , we see that the square of the determinant of the Jacobian matrix is

$$|\mathbf{J}|^2 = |\mathbf{I}| = 1 \quad (4.21)$$

Also, the determinant  $|\mathbf{\Sigma}|$  of the covariance matrix can be written as the product of its eigenvalues, and hence

$$|\Sigma|^{1/2} = \prod_{j=1}^D \lambda_j^{1/2} \quad (4.22)$$

Thus in the  $y_i$  coordinate system, the Gaussian distribution takes the form

$$p(\mathbf{y}) = p(\mathbf{x})|\mathbf{J}| = \prod_{j=1}^D \frac{1}{(2\pi\lambda_j)^{1/2}} \exp\left\{-\frac{y_j^2}{2\lambda_j}\right\} \quad (4.23)$$

which is the product of D independent univariate Gaussian distributions. So we have normalized the multivariate Gaussian.

Now we can check the first and second order moments of the Gaussian.

## 4.1.2 limitations

Gaussian distribution has quadratically growing parameters with dimension. A further limitation of the Gaussian distribution is that it is intrinsically unimodal (i.e., has a single maximum) and so is unable to provide a good approximation to multimodal distributions. We will introduce *latent* variables, also called *hidden* variables or *unobserved* variables to address both of these problems.

### 4.1.3 MLE for a MVN

**Theorem 4.1. (MLE for a MVN)** *If we have  $N$  iid samples  $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then the MLE for the parameters is given by*

$$\bar{\boldsymbol{\mu}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \triangleq \bar{\mathbf{x}} \quad (4.24)$$

$$\bar{\boldsymbol{\Sigma}} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (4.25)$$

$$= \frac{1}{N} \left( \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) - \bar{\mathbf{x}} \bar{\mathbf{x}}^T \quad (4.26)$$

### 4.1.4 Maximum entropy derivation of the Gaussian \*

In this section, we show that the multivariate Gaussian is the distribution with maximum entropy subject to having a specified mean and covariance (see also Section TODO). This is one reason the Gaussian is so widely used: the first two moments are usually all that we can reliably estimate from data, so we want a distribution that captures these properties, but otherwise makes as few additional assumptions as possible.

To simplify notation, we will assume the mean is zero. The pdf has the form

$$f(\mathbf{x}) = \frac{1}{Z} \exp \left( -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} \right) \quad (4.27)$$

## 4.2 Gaussian discriminant analysis

One important application of MVNs is to define the the class conditional densities in a generative classifier, i.e.,

$$p(\mathbf{x} | y = c, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \quad (4.28)$$

The resulting technique is called (Gaussian) **discriminant analysis** or **GDA** (even though it is a generative, not discriminative, classifier see Section TODO for more on this distinction). If  $\boldsymbol{\Sigma}_c$  is diagonal, this is equivalent to naive Bayes.

We can classify a feature vector using the following decision rule, derived from Equation 3.1:

$$y = \arg \max_c [\log p(y = c | \boldsymbol{\pi}) + \log p(\mathbf{x} | \boldsymbol{\theta})] \quad (4.29)$$

When we compute the probability of  $\mathbf{x}$  under each class conditional density, we are measuring the distance from  $\mathbf{x}$  to the center of each class,  $\boldsymbol{\mu}_c$ , using Mahalanobis distance. This can be thought of as a **nearest centroids classifier**.

As an example, Figure 4.2 shows two Gaussian class-conditional densities in 2d, representing the height and weight of men and women. We can see that the features are correlated, as is to be expected (tall people tend to weigh more). The ellipses for each class contain 95% of the probability mass. If we have a uniform prior over classes, we can classify a new test vector as follows:

$$y = \arg \max_c (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \quad (4.30)$$

### 4.2.1 Quadratic discriminant analysis (QDA)

By plugging in the definition of the Gaussian density to Equation 3.1, we can get



Fig. 4.2: (a) Height/weight data. (b) Visualization of 2d Gaussians fit to each class. 95% of the probability mass is inside the ellipse.

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \frac{\pi_c |2\pi \boldsymbol{\Sigma}_c|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]}{\sum_{c'} \pi_{c'} |2\pi \boldsymbol{\Sigma}_{c'}|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]} \quad (4.31)$$

Thresholding this results in a quadratic function of  $\mathbf{x}$ . The result is known as quadratic discriminant analysis (QDA). Figure 4.3 gives some examples of what the decision boundaries look like in 2D.

## 4.2.2 Linear discriminant analysis (LDA)

We now consider a special case in which the covariance matrices are **tied** or **shared** across classes,  $\boldsymbol{\Sigma}_c = \boldsymbol{\Sigma}$ . In this case, we can simplify Equation 4.31 as follows:



Fig. 4.3: Quadratic decision boundaries in 2D for the 2 and 3 class case.

$$\begin{aligned}
 p(y|\mathbf{x}, \boldsymbol{\theta}) &\propto \pi_c \exp \left( \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c \right) \\
 &= \exp \left( \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \right) \\
 &\quad \exp \left( -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} \right) \\
 &\propto \exp \left( \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \right)
 \end{aligned} \tag{4.32}$$

Since the quadratic term  $\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}$  is independent of  $c$ , it will cancel out in the numerator and denominator. If we define

$$\gamma_c \triangleq -\frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \tag{4.33}$$

$$\boldsymbol{\beta}_c \triangleq \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c \tag{4.34}$$

then we can write

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\boldsymbol{\beta}_c^T \mathbf{x} + \gamma_c}}{\sum_{c'} e^{\boldsymbol{\beta}_{c'}^T \mathbf{x} + \gamma_{c'}}} \triangleq \sigma(\boldsymbol{\eta}, c) \tag{4.35}$$



where  $\boldsymbol{\eta} \triangleq (e^{\beta_1^T \mathbf{x}} + \gamma_1, \dots, e^{\beta_C^T \mathbf{x}} + \gamma_C)$ ,  $\sigma(\cdot)$  is the **softmax activation function**<sup>17</sup>, defined as follows:

$$\sigma(\mathbf{q}, i) \triangleq \frac{\exp(q_i)}{\sum_{j=1}^n \exp(q_j)} \quad (4.36)$$

When parameterized by some constant,  $\alpha > 0$ , the following formulation becomes a smooth, differentiable approximation of the maximum function:

$$\mathcal{S}_\alpha(\mathbf{x}) = \frac{\sum_{j=1}^D x_j e^{\alpha x_j}}{\sum_{j=1}^D e^{\alpha x_j}} \quad (4.37)$$

$\mathcal{S}_\alpha$  has the following properties:

1.  $\mathcal{S}_\alpha \rightarrow \max$  as  $\alpha \rightarrow \infty$
2.  $\mathcal{S}_0$  is the average of its inputs
3.  $\mathcal{S}_\alpha \rightarrow \min$  as  $\alpha \rightarrow -\infty$

Note that the softmax activation function comes from the area of statistical physics, where it is common to use the **Boltzmann distribution**, which has the same form as the softmax activation function.

An interesting property of Equation 4.35 is that, if we take logs, we end up with a linear function of  $\mathbf{x}$ . (The reason it is linear is because the  $\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}$  cancels from the numerator and denominator.) Thus the decision boundary between any two classes, says  $c$  and  $c'$ , will be a straight line. Hence this technique is called **linear discriminant analysis** or **LDA**.

An alternative to fitting an LDA model and then deriving the class posterior is to directly fit  $p(y|\mathbf{x}, \mathbf{W}) = \text{Cat}(y|\mathbf{W}\mathbf{x})$  for some  $C \times D$  weight matrix  $\mathbf{W}$ . This is called **multi-class logistic regression**, or **multinomial logistic regression**. We will discuss this model in detail in Section TODO. The difference between the two approaches is explained in Section TODO.

### 4.2.3 Two-class LDA

To gain further insight into the meaning of these equations, let us consider the binary case. In this case, the posterior is given by

$$p(y = 1|\mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\beta_1^T \mathbf{x} + \gamma_1}}{e^{\beta_0^T \mathbf{x} + \gamma_0} + e^{\beta_1^T \mathbf{x} + \gamma_1}} \quad (4.38)$$

$$= \frac{1}{1 + e^{(\beta_0 - \beta_1)^T \mathbf{x} + (\gamma_0 - \gamma_1)}} \quad (4.39)$$

$$= \text{sigm}((\beta_1 - \beta_0)^T \mathbf{x} + (\gamma_0 - \gamma_1)) \quad (4.40)$$

where  $\text{sigm}(x)$  refers to the sigmoid function<sup>18</sup>.

Now

$$\gamma_1 - \gamma_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 + \log(\pi_1/\pi_0) \quad (4.41)$$

$$= -\frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) + \log(\pi_1/\pi_0) \quad (4.42)$$

So if we define

$$\mathbf{w} = \beta_1 - \beta_0 = \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \quad (4.43)$$

$$\mathbf{x}_0 = \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) - (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \frac{\log(\pi_1/\pi_0)}{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)} \quad (4.44)$$

<sup>17</sup> [http://en.wikipedia.org/wiki/Softmax\\_activation\\_function](http://en.wikipedia.org/wiki/Softmax_activation_function)

<sup>18</sup> [http://en.wikipedia.org/wiki/Sigmoid\\_function](http://en.wikipedia.org/wiki/Sigmoid_function)

then we have  $\mathbf{w}^T \mathbf{x}_0 = -(\gamma_1 - \gamma_0)$ , and hence

$$p(y = 1 | \mathbf{x}, \boldsymbol{\theta}) = \text{sigm}(\mathbf{w}^T (\mathbf{x} - \mathbf{x}_0)) \quad (4.45)$$

(This is closely related to logistic regression, which we will discuss in Section TODO.) So the final decision rule is as follows: shift  $\mathbf{x}$  by  $\mathbf{x}_0$ , project onto the line  $\mathbf{w}$ , and see if the result is positive or negative.



Fig. 4.4: Geometry of LDA in the 2 class case where  $\Sigma_1 = \Sigma_2 = I$ .

If  $\Sigma = \sigma^2 I$ , then  $\mathbf{w}$  is in the direction of  $\mu_1 - \mu_0$ . So we classify the point based on whether its projection is closer to  $\mu_0$  or  $\mu_1$ . This is illustrated in Figure 4.4. Furthermore, if  $\pi_1 = \pi_0$ , then  $\mathbf{x}_0 = \frac{1}{2}(\mu_1 + \mu_0)$ , which is half way between the means. If we make  $\pi_1 > \pi_0$ , then  $\mathbf{x}_0$  gets closer to  $\mu_0$ , so more of the line belongs to class 1 *a priori*. Conversely if  $\pi_1 < \pi_0$ , the boundary shifts right. Thus we see that the class prior,  $c$ , just changes the decision threshold, and not the overall geometry, as we claimed above. (A similar argument applies in the multi-class case.)

The magnitude of  $\mathbf{w}$  determines the steepness of the logistic function, and depends on how well-separated the means are, relative to the variance. In psychology and signal detection theory, it is common to define the **discriminability** of a signal from the background noise using a quantity called **d-prime**:

$$d' \triangleq \frac{\mu_1 - \mu_0}{\sigma} \quad (4.46)$$

where  $\mu_1$  is the mean of the signal and  $\mu_0$  is the mean of the noise, and  $\sigma$  is the standard deviation of the noise. If  $d'$  is large, the signal will be easier to discriminate from the noise.

#### 4.2.4 MLE for discriminant analysis

The log-likelihood function is as follows:

$$p(\mathcal{D} | \boldsymbol{\theta}) = \sum_{c=1}^C \sum_{i: y_i=c} \log \pi_c + \sum_{c=1}^C \sum_{i: y_i=c} \log \mathcal{N}(\mathbf{x}_i | \mu_c, \Sigma_c) \quad (4.47)$$

The MLE for each parameter is as follows:

$$\bar{\mu}_c = \frac{N_c}{N} \quad (4.48)$$

$$\bar{\mu}_c = \frac{1}{N_c} \sum_{i: y_i=c} \mathbf{x}_i \quad (4.49)$$

$$\bar{\Sigma}_c = \frac{1}{N_c} \sum_{i: y_i=c} (\mathbf{x}_i - \bar{\mu}_c)(\mathbf{x}_i - \bar{\mu}_c)^T \quad (4.50)$$

### 4.2.5 Strategies for preventing overfitting

The speed and simplicity of the MLE method is one of its greatest appeals. However, the MLE can badly overfit in high dimensions. In particular, the MLE for a full covariance matrix is singular if  $N_c < D$ . And even when  $N_c > D$ , the MLE can be ill-conditioned, meaning it is close to singular. There are several possible solutions to this problem:

- Use a diagonal covariance matrix for each class, which assumes the features are conditionally independent; this is equivalent to using a naive Bayes classifier (Section 3.5).
- Use a full covariance matrix, but force it to be the same for all classes,  $\Sigma_c = \Sigma$ . This is an example of **parameter tying** or **parameter sharing**, and is equivalent to LDA (Section 4.2.2).
- Use a diagonal covariance matrix and forced it to be shared. This is called diagonal covariance LDA, and is discussed in Section TODO.
- Use a full covariance matrix, but impose a prior and then integrate it out. If we use a conjugate prior, this can be done in closed form, using the results from Section TODO; this is analogous to the Bayesian naive Bayes method in Section 3.5.1.2. See (Minka 2000f) for details.
- Fit a full or diagonal covariance matrix by MAP estimation. We discuss two different kindsof prior below.
- Project the data into a low dimensional subspace and fit the Gaussians there. See Section TODO for a way to find the best (most discriminative) linear projection.

We discuss some of these options below.

### 4.2.6 Regularized LDA \*

### 4.2.7 Diagonal LDA

### 4.2.8 Nearest shrunken centroids classifier \*

One drawback of diagonal LDA is that it depends on all of the features. In high dimensional problems, we might prefer a method that only depends on a subset of the features, for reasons of accuracy and interpretability. One approach is to use a screening method, perhaps based on mutual information, as in Section 3.5.4. We now discuss another approach to this problem known as the **nearest shrunken centroids** classifier (Hastie et al. 2009, p652).

## 4.3 Inference in jointly Gaussian distributions

Given a joint distribution,  $p(\mathbf{x}_1, \mathbf{x}_2)$ , it is useful to be able to compute marginals  $p(\mathbf{x}_1)$  and conditionals  $p(\mathbf{x}_1|\mathbf{x}_2)$ . We discuss how to do this below, and then give some applications. These operations take  $O(D^3)$  time in the worst case. See Section TODO for faster methods.

### 4.3.1 Statement of the result

**Theorem 4.2. (Marginals and conditionals of an MVN).** Suppose  $X = (\mathbf{x}_1, \mathbf{x}_2)$  is jointly Gaussian with parameters

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}, \boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1} = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix}, \quad (4.51)$$

Then the marginals are given by

$$\begin{aligned} p(\mathbf{x}_1) &= \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}) \\ p(\mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_2 | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}) \end{aligned} \tag{4.52}$$

and the posterior conditional is given by

$$\begin{aligned} p(\mathbf{x}_1 | \mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2}) \\ \boldsymbol{\mu}_{1|2} &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{11}^{-1} \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\Sigma}_{1|2} (\boldsymbol{\Lambda}_{11} \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2)) \\ \boldsymbol{\Sigma}_{1|2} &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} = \boldsymbol{\Lambda}_{11}^{-1} \end{aligned} \tag{4.53}$$

Equation 4.53 is of such crucial importance in this book that we have put a box around it, so you can easily find it. For the proof, see Section TODO.

We see that both the marginal and conditional distributions are themselves Gaussian. For the marginals, we just extract the rows and columns corresponding to  $\mathbf{x}_1$  or  $\mathbf{x}_2$ . For the conditional, we have to do a bit more work. However, it is not that complicated: the conditional mean is just a linear function of  $\mathbf{x}_2$ , and the conditional covariance is just a constant matrix that is independent of  $\mathbf{x}_2$ . We give three different (but equivalent) expressions for the posterior mean, and two different (but equivalent) expressions for the posterior covariance; each one is useful in different circumstances.

## 4.3.2 Examples

Below we give some examples of these equations in action, which will make them seem more intuitive.

### 4.3.2.1 Marginals and conditionals of a 2d Gaussian

### 4.3.3 information form

### 4.3.4 Proof

## 4.4 Linear Gaussian systems

Suppose we have two variables,  $\mathbf{x}$  and  $\mathbf{y}$ . Let  $\mathbf{x} \in \mathbb{R}^{D_x}$  be a hidden variable, and  $\mathbf{y} \in \mathbb{R}^{D_y}$  be a noisy observation of  $\mathbf{x}$ . Let us assume we have the following prior and likelihood:

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x) \\ p(\mathbf{y} | \mathbf{x}) &= \mathcal{N}(\mathbf{y} | \mathbf{W}\mathbf{x} + \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) \end{aligned} \tag{4.54}$$

where  $\mathbf{W}$  is a matrix of size  $D_y \times D_x$ . This is an example of a **linear Gaussian system**. We can represent this schematically as  $\mathbf{x} \rightarrow \mathbf{y}$ , meaning  $\mathbf{x}$  generates  $\mathbf{y}$ . In this section, we show how to invert the arrow, that is, how to infer  $\mathbf{x}$  from  $\mathbf{y}$ . We state the result below, then give several examples, and finally we derive the result. We will see many more applications of these results in later chapters.

### 4.4.1 Statement of the result

**Theorem 4.3. (Bayes rule for linear Gaussian systems).** Given a linear Gaussian system, as in Equation 4.54, the posterior  $p(\mathbf{x}|\mathbf{y})$  is given by the following:

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y}) \\ \boldsymbol{\Sigma}_{x|y} &= \boldsymbol{\Sigma}_x^{-1} + \mathbf{W}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{W} \\ \boldsymbol{\mu}_{x|y} &= \boldsymbol{\Sigma}_{x|y} [\mathbf{W}^T \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y) + \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\mu}_x] \end{aligned} \quad (4.55)$$

In addition, the normalization constant  $p(\mathbf{y})$  is given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{W}\boldsymbol{\mu}_x + \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y + \mathbf{W}\boldsymbol{\Sigma}_x\mathbf{W}^T) \quad (4.56)$$

For the proof, see Section 4.4.3 TODO.

### 4.4.2 Examples

### 4.4.3 Proof

### 4.5 Digression: The Wishart distribution \*

### 4.6 Inferring the parameters of an MVN

#### 4.6.1 Posterior distribution of $\boldsymbol{\mu}$

#### 4.6.2 Posterior distribution of $\boldsymbol{\Sigma}$ \*

#### 4.6.3 Posterior distribution of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ \*

#### 4.6.4 Sensor fusion with unknown precisions \*



## Chapter 5

# Bayesian statistics

### 5.1 Introduction

Using the posterior distribution to summarize everything we know about a set of unknown variables is at the core of Bayesian statistics. In this chapter, we discuss this approach to statistics in more detail.

### 5.2 Summarizing posterior distributions

The posterior  $p(\theta|\mathcal{D})$  summarizes everything we know about the unknown quantities  $\theta$ . In this section, we discuss some simple quantities that can be derived from a probability distribution, such as a posterior. These summary statistics are often easier to understand and visualize than the full joint.

#### 5.2.1 MAP estimation

We can easily compute a **point estimate** of an unknown quantity by computing the posterior mean, median or mode. In Section 5.7, we discuss how to use decision theory to choose between these methods. Typically the posterior mean or median is the most appropriate choice for a realvalued quantity, and the vector of posterior marginals is the best choice for a discrete quantity. However, the posterior mode, aka the MAP estimate, is the most popular choice because it reduces to an optimization problem, for which efficient algorithms often exist. Furthermore, MAP estimation can be interpreted in non-Bayesian terms, by thinking of the log prior as a regularizer (see Section TODO for more details).

Although this approach is computationally appealing, it is important to point out that there are various drawbacks to MAP estimation, which we briefly discuss below. This will provide motivation for the more thoroughly Bayesian approach which we will study later in this chapter(and elsewhere in this book).

##### 5.2.1.1 No measure of uncertainty

The most obvious drawback of MAP estimation, and indeed of any other *point estimate* such as the posterior mean or median, is that it does not provide any measure of uncertainty. In many applications, it is important to know how much one can trust a given estimate. We can derive such confidence measures from the posterior, as we discuss in Section 5.2.2.

##### 5.2.1.2 Plugging in the MAP estimate can result in overfitting

If we don't model the uncertainty in our parameters, then our predictive distribution will be overconfident. Overconfidence in predictions is particularly problematic in situations where we may be risk averse; see Section 5.7 for details.

##### 5.2.1.3 The mode is an untypical point

Choosing the mode as a summary of a posterior distribution is often a very poor choice, since the mode is usually quite untypical of the distribution, unlike the mean or median. The basic problem is that the mode is a point of measure zero, whereas the mean and median take the volume of the space into account. See Figure 5.1.



Fig. 5.1: (a) A bimodal distribution in which the mode is very untypical of the distribution. The thin blue vertical line is the mean, which is arguably a better summary of the distribution, since it is near the majority of the probability mass. (b) A skewed distribution in which the mode is quite different from the mean.

How should we summarize a posterior if the mode is not a good choice? The answer is to use decision theory, which we discuss in Section 5.7. The basic idea is to specify a loss function, where  $L(\theta, \hat{\theta})$  is the loss you incur if the truth is  $\theta$  and your estimate is  $\hat{\theta}$ . If we use 0-1 loss  $L(\theta, \hat{\theta}) = \mathbb{I}(\theta \neq \hat{\theta})$  (see section 1.4.2.1), then the optimal estimate is the posterior mode. 0-1 loss means you only get points if you make no errors, otherwise you get nothing: there is no partial credit under this loss function! For continuous-valued quantities, we often prefer to use squared error loss,  $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$ ; the corresponding optimal estimator is then the posterior mean, as we show in Section 5.7. Or we can use a more robust loss function,  $L(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$ , which gives rise to the posterior median.

#### 5.2.1.4 MAP estimation is not invariant to reparameterization \*

A more subtle problem with MAP estimation is that the result we get depends on how we parameterize the probability distribution. Changing from one representation to another equivalent representation changes the result, which is not very desirable, since the units of measurement are arbitrary (e.g., when measuring distance, we can use centimetres or inches).

To understand the problem, suppose we compute the posterior for  $x$ . If we define  $y=f(x)$ , the distribution for  $y$  is given by Equation 2.100. The  $\frac{dx}{dy}$  term is called the Jacobian, and it measures the change in size of a unit volume passed through  $f$ . Let  $\hat{x} = \arg \max_x p_x(x)$  be the MAP estimate for  $x$ . In general it is not the case that  $\hat{x} = \arg \max_x p_x(x)$  is given by  $f(\hat{x})$ . For example, let  $X \sim \mathcal{N}(6, 1)$  and  $y = f(x)$ , where  $f(x) = 1/(1 + \exp(-x + 5))$ .

We can derive the distribution of  $y$  using Monte Carlo simulation (see Section 2.10). The result is shown in Figure ???. We see that the original Gaussian has become squashed by the sigmoid nonlinearity. In particular, we see that the mode of the transformed distribution is not equal to the transform of the original mode.

The MLE does not suffer from this since the likelihood is a function, not a probability density. Bayesian inference does not suffer from this problem either, since the change of measure is taken into account when integrating over the parameter space.





Fig. 5.2: Example of the transformation of a density under a nonlinear transform. Note how the mode of the transformed distribution is not the transform of the original mode. Based on Exercise 1.4 of (Bishop 2006b).

## 5.2.2 Credible intervals

In addition to point estimates, we often want a measure of confidence. A standard measure of confidence in some (scalar) quantity  $\theta$  is the width of its posterior distribution. This can be measured using a  $100(1-\alpha)\%$  credible interval, which is a (contiguous) region  $C = (\ell, u)$  (standing for lower and upper) which contains  $1-\alpha$  of the posterior probability mass, i.e.,

$$C_\alpha(\mathcal{D}) \quad \text{where } P(\ell \leq \theta \leq u) = 1 - \alpha \quad (5.1)$$

There may be many such intervals, so we choose one such that there is  $(1-\alpha)/2$  mass in each tail; this is called a **central interval**.

If the posterior has a known functional form, we can compute the posterior central interval using  $\ell = F^{-1}(\alpha/2)$  and  $u = F^{-1}(1 - \alpha/2)$ , where  $F$  is the cdf of the posterior.

If we don't know the functional form, but we can draw samples from the posterior, then we can use a Monte Carlo approximation to the posterior quantiles: we simply sort the  $S$  samples, and find the one that occurs at location  $\alpha/S$  along the sorted list. As  $S \rightarrow \infty$ , this converges to the true quantile.

People often confuse Bayesian credible intervals with frequentist confidence intervals. However, they are not the same thing, as we discuss in Section TODO. In general, credible intervals are usually what people want to compute, but confidence intervals are usually what they actually compute, because most people are taught frequentist statistics but not Bayesian statistics. Fortunately, the mechanics of computing a credible interval is just as easy as computing a confidence interval.

## 5.2.3 Inference for a difference in proportions

Sometimes we have multiple parameters, and we are interested in computing the posterior distribution of some function of these parameters. For example, suppose you are about to buy something from Amazon.com, and there are two sellers offering it for the same price. Seller 1 has 90 positive reviews and 10 negative reviews. Seller 2 has 2 positive reviews and 0 negative reviews. Who should you buy from?<sup>19</sup>

On the face of it, you should pick seller 2, but we cannot be very confident that seller 2 is better since it has had so few reviews. In this section, we sketch a Bayesian analysis of this problem. Similar methodology can be used to compare rates or proportions across groups for a variety of other settings.

<sup>19</sup> This example is from <http://www.johndcook.com/blog/2011/09/27/bayesian-amazon/>

Let  $\theta_1$  and  $\theta_2$  be the unknown reliabilities of the two sellers. Since we don't know much about them, we'll endow them both with uniform priors,  $\theta_i \sim \text{Beta}(1, 1)$ . The posteriors are  $p(\theta_1|\mathcal{D}_1) = \text{Beta}(91, 11)$  and  $p(\theta_2|\mathcal{D}_2) = \text{Beta}(3, 1)$ .

We want to compute  $p(\theta_1 > \theta_2|\mathcal{D})$ . For convenience, let us define  $\delta = \theta_1 - \theta_2$  as the difference in the rates. (Alternatively we might want to work in terms of the log-odds ratio.) We can compute the desired quantity using numerical integration

$$p(\delta > 0|\mathcal{D}) = \int_0^1 \int_0^1 \mathbb{I}(\theta_1 > \theta_2) \text{Beta}(\theta_1|91, 11) \text{Beta}(\theta_2|3, 1) d\theta_1 d\theta_2 \quad (5.2)$$

We find  $p(\delta > 0|\mathcal{D}) = 0.710$ , which means you are better off buying from seller 1!

### 5.3 Bayesian model selection

In general, when faced with a set of models (i.e., families of parametric distributions) of different complexity, how should we choose the best one? This is called the **model selection** problem.

One approach is to use cross-validation to estimate the generalization error of all the candidate models, and then to pick the model that seems the best. However, this requires fitting each model  $K$  times, where  $K$  is the number of CV folds. A more efficient approach is to compute the posterior over models,

$$p(m|\mathcal{D}) = \frac{p(\mathcal{D}|m)p(m)}{\sum_{m'} p(\mathcal{D}|m')p(m')} \quad (5.3)$$

From this, we can easily compute the MAP model,  $\hat{m} = \arg \max_m p(m|\mathcal{D})$ . This is called **Bayesian model selection**. If we use a uniform prior over models, this amounts to picking the model which maximizes

$$p(\mathcal{D}|m) = \int p(\mathcal{D}|\theta)p(\theta|m)d\theta \quad (5.4)$$

This quantity is called the **marginal likelihood**, the **integrated likelihood**, or the **evidence** for model  $m$ . The details on how to perform this integral will be discussed in Section 5.3.2. But first we give an intuitive interpretation of what this quantity means.

#### 5.3.1 Bayesian Occam's razor

One might think that using  $p(\mathcal{D}|m)$  to select models would always favour the model with the most parameters. This is true if we use  $p(\mathcal{D}|\hat{\theta}_m)$  to select models, where  $\hat{\theta}_m$  is the MLE or MAP estimate of the parameters for model  $m$ , because models with more parameters will fit the data better, and hence achieve higher likelihood. However, if we integrate out the parameters, rather than maximizing them, we are automatically protected from overfitting: models with more parameters do not necessarily have higher *marginal likelihood*. This is called the **Bayesian Occams razor** effect (MacKay 1995b; Murray and Ghahramani 2005), named after the principle known as **Occams razor**, which says one should pick the simplest model that adequately explains the data.

One way to understand the Bayesian Occams razor is to notice that the marginal likelihood can be rewritten as follows, based on the chain rule of probability (Equation 2.13):

$$\begin{aligned} p(D) &= p((x_1, y_1))p((x_2, y_2)|(x_1, y_1)) \\ &\quad p((x_3, y_3)|(x_1, y_1) : (x_2, y_2)) \cdots \\ &\quad p((x_N, y_N)|(x_1, y_1) : (x_{N-1}, y_{N-1})) \end{aligned} \quad (5.5)$$

This is similar to a leave-one-out cross-validation estimate (Section 1.5.0.1) of the likelihood, since we predict each future point given all the previous ones. (Of course, the order of the data does not matter in the above expression.) If a model is too complex, it will overfit the early examples and will then predict the remaining ones poorly.

Another way to understand the Bayesian Occams razor effect is to note that probabilities must sum to one. Hence  $\sum_{p(\mathcal{D}')} p(m|\mathcal{D}') = 1$ , where the sum is over all possible data sets. Complex models, which can predict many things, must spread their probability mass thinly, and hence will not obtain as large a probability for any given data set as simpler models. This is sometimes called the **conservation of probability mass** principle, and is illustrated in Figure 5.3.

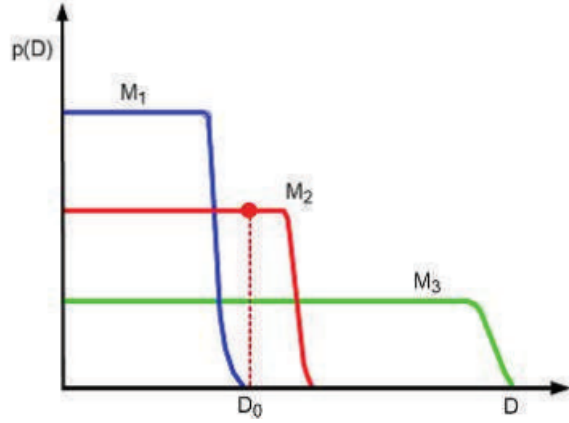


Fig. 5.3: A schematic illustration of the Bayesian Occams razor. The broad (green) curve corresponds to a complex model, the narrow (blue) curve to a simple model, and the middle (red) curve is just right. Based on Figure 3.13 of (Bishop 2006a).

When using the Bayesian approach, we are not restricted to evaluating the evidence at a finite grid of values. Instead, we can use numerical optimization to find  $\lambda^* = \arg \max_{\lambda} p(\mathcal{D}|\lambda)$ . This technique is called **empirical Bayes** or **type II maximum likelihood** (see Section 5.6 for details). An example is shown in Figure TODO(b): we see that the curve has a similar shape to the CV estimate, but it can be computed more efficiently.

### 5.3.2 Computing the marginal likelihood (evidence)

When discussing parameter inference for a fixed model, we often wrote

$$p(\theta|\mathcal{D}, m) \propto p(\theta|m)p(\mathcal{D}|\theta, m) \quad (5.6)$$

thus ignoring the normalization constant  $p(\mathcal{D}|m)$ . This is valid since  $p(\mathcal{D}|m)$  is constant wrt  $\theta$ . However, when comparing models, we need to know how to compute the marginal likelihood,  $p(\mathcal{D}|m)$ . In general, this can be quite hard, since we have to integrate over all possible parameter values, but when we have a conjugate prior, it is easy to compute, as we now show.

Let  $p(\theta) = q(\theta)/Z_0$  be our prior, where  $q(\theta)$  is an unnormalized distribution, and  $Z_0$  is the normalization constant of the prior. Let  $p(\mathcal{D}|\theta) = q(\mathcal{D}|\theta)/Z_\ell$  be the likelihood, where  $Z_\ell$  contains any constant factors in the likelihood. Finally let  $p(\theta|\mathcal{D}) = q(\theta|\mathcal{D})/Z_N$  be our posterior, where  $q(\theta|\mathcal{D}) = q(\mathcal{D}|\theta)q(\theta)$  is the unnormalized posterior, and  $Z_N$  is the normalization constant of the posterior. We have

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})} \quad (5.7)$$

$$\frac{q(\boldsymbol{\theta}|\mathcal{D})}{Z_N} = \frac{q(\mathcal{D}|\boldsymbol{\theta})q(\boldsymbol{\theta})}{Z_\ell Z_0 p(\mathcal{D})} \quad (5.8)$$

$$p(\mathcal{D}) = \frac{Z_N}{Z_0 Z_\ell} \quad (5.9)$$

So assuming the relevant normalization constants are tractable, we have an easy way to compute the marginal likelihood. We give some examples below.

### 5.3.2.1 Beta-binomial model

Let us apply the above result to the Beta-binomial model. Since we know  $p(\boldsymbol{\theta}|\mathcal{D}) = \text{Beta}(\boldsymbol{\theta}|a', b')$ , where  $a' = a + N_1$ ,  $b' = b + N_0$ , we know the normalization constant of the posterior is  $B(a', b')$ . Hence

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})} \quad (5.10)$$

$$= \frac{1}{p(\mathcal{D})} \left[ \frac{1}{B(a, b)} \theta^{a-1} (1-\theta)^{b-1} \right] \left[ \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N_0} \right] \quad (5.11)$$

$$= \binom{N}{N_1} \frac{1}{p(\mathcal{D})} \frac{1}{B(a, b)} \left[ \theta^{a+N_1-1} (1-\theta)^{b+N_0-1} \right] \quad (5.12)$$

So

$$\frac{1}{B(a+N_1, b+N_0)} = \binom{N}{N_1} \frac{1}{p(\mathcal{D})} \frac{1}{B(a, b)} \quad (5.13)$$

$$p(\mathcal{D}) = \binom{N}{N_1} \frac{B(a+N_1, b+N_0)}{B(a, b)} \quad (5.14)$$

The marginal likelihood for the Beta-Bernoulli model is the same as above, except it is missing the  $\binom{N}{N_1}$  term.

### 5.3.2.2 Dirichlet-multinoulli model

By the same reasoning as the Beta-Bernoulli case, one can show that the marginal likelihood for the Dirichlet-multinoulli model is given by

$$p(\mathcal{D}) = \frac{B(\mathbf{N} + \boldsymbol{\alpha})}{B(\boldsymbol{\alpha})} \quad (5.15)$$

$$= \frac{\Gamma(\sum_k \alpha_k)}{\Gamma(N + \sum_k \alpha_k)} \prod_k \frac{\Gamma(N_k + \alpha_k)}{\Gamma(\alpha_k)} \quad (5.16)$$

### 5.3.2.3 Gaussian-Gaussian-Wishart model

Consider the case of an MVN with a conjugate NIW prior. Let  $Z_0$  be the normalizer for the prior,  $Z_N$  be normalizer for the posterior, and let  $Z_\ell(2\pi)^{ND/2}$  be the normalizer for the likelihood. Then it is easy to see that

$$p(\mathcal{D}) = \frac{Z_N}{Z_0 Z_\ell} \quad (5.17)$$

$$= \frac{1}{(2\pi)^{ND/2}} \frac{\left(\frac{2\pi}{\kappa_N}\right)^{D/2} |\mathbf{S}_N|^{-v_N/2} 2^{(v_0+N)D/2} \Gamma_D(v_N/2)}{\left(\frac{2\pi}{\kappa_0}\right)^{D/2} |\mathbf{S}_0|^{-v_0/2} 2^{v_0 D/2} \Gamma_D(v_0/2)} \quad (5.18)$$

$$= \frac{1}{\pi^{ND/2}} \left(\frac{\kappa_0}{\kappa_N}\right)^{D/2} \frac{|\mathbf{S}_0|^{v_0/2} \Gamma_D(v_N/2)}{|\mathbf{S}_N|^{v_N/2} \Gamma_D(v_0/2)} \quad (5.19)$$

### 5.3.2.4 BIC approximation to log marginal likelihood

In general, computing the integral in Equation 5.4 can be quite difficult. One simple but popular approximation is known as the **Bayesian information criterion** or **BIC**, which has the following form (Schwarz 1978):

$$\text{BIC} \triangleq \log p(\mathcal{D}|\hat{\boldsymbol{\theta}}) - \frac{\text{dof}(\hat{\boldsymbol{\theta}})}{2} \log N \quad (5.20)$$

where  $\text{dof}(\hat{\boldsymbol{\theta}})$  is the number of **degrees of freedom** in the model, and  $\hat{\boldsymbol{\theta}}$  is the MLE for the model. We see that this has the form of a **penalized log likelihood**, where the penalty term depends on the models complexity. See Section TODO for the derivation of the BIC score.

As an example, consider linear regression. As we show in Section TODO, the MLE is given by  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  and  $\sigma^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\mathbf{w}}^T \mathbf{x}_i)^2$ . The corresponding log likelihood is given by

$$\log p(\mathcal{D}|\hat{\boldsymbol{\theta}}) = -\frac{N}{2} \log(2\pi\hat{\sigma}^2) - \frac{N}{2} \quad (5.21)$$

Hence the BIC score is as follows (dropping constant terms)

$$\text{BIC} = -\frac{N}{2} \log(\hat{\sigma}^2) - \frac{D}{2} \log N \quad (5.22)$$

where  $D$  is the number of variables in the model. In the statistics literature, it is common to use an alternative definition of BIC, which we call the **BIC cost** (since we want to minimize it):

$$\text{BIC-cost} \triangleq -2 \log p(\mathcal{D}|\hat{\boldsymbol{\theta}}) - \text{dof}(\hat{\boldsymbol{\theta}}) \log N \approx -2 \log p(\mathcal{D}) \quad (5.23)$$

In the context of linear regression, this becomes

$$\text{BIC-cost} = N \log(\hat{\sigma}^2) + D \log N \quad (5.24)$$

The BIC method is very closely related to the **minimum description length** or **MDL** principle, which characterizes the score for a model in terms of how well it fits the data, minus how complex the model is to define. See (Hansen and Yu 2001) for details.

There is a very similar expression to BIC/ MDL called the **Akaike information criterion** or **AIC**, defined as

$$\text{AIC}(m, \mathcal{D}) = \log p(\mathcal{D}|\hat{\boldsymbol{\theta}}_{MLE}) - \text{dof}(m) \quad (5.25)$$

This is derived from a frequentist framework, and cannot be interpreted as an approximation to the marginal likelihood. Nevertheless, the form of this expression is very similar to BIC. We see that the penalty for AIC is less than for BIC. This causes AIC to pick more complex models. However, this can result in better predictive accuracy. See e.g., (Clarke et al. 2009, sec 10.2) for further discussion on such information criteria.

### 5.3.2.5 Effect of the prior

Sometimes it is not clear how to set the prior. When we are performing posterior inference, the details of the prior may not matter too much, since the likelihood often overwhelms the prior anyway. But when computing the marginal likelihood, the prior plays a much more important role, since we are averaging the likelihood over all possible parameter settings, as weighted by the prior.

If the prior is unknown, the correct Bayesian procedure is to put a prior on the prior. If the prior is unknown, the correct Bayesian procedure is to put a prior on the prior.

## 5.3.3 Bayes factors

Suppose our prior on models is uniform,  $p(m) \propto 1$ . Then model selection is equivalent to picking the model with the highest marginal likelihood. Now suppose we just have two models we are considering, call them the **null hypothesis**,  $M_0$ , and the **alternative hypothesis**,  $M_1$ . Define the **Bayes factor** as the ratio of marginal likelihoods:

$$\text{BF}_{1,0} \triangleq \frac{p(\mathcal{D}|M_1)}{p(\mathcal{D}|M_0)} = \frac{p(M_1|\mathcal{D})}{p(M_2|\mathcal{D})} \frac{p(M_1)}{p(M_0)} \quad (5.26)$$

## 5.4 Priors

The most controversial aspect of Bayesian statistics is its reliance on priors. Bayesians argue this is unavoidable, since nobody is a **tabula rasa** or **blank slate**: all inference must be done conditional on certain assumptions about the world. Nevertheless, one might be interested in minimizing the impact of one's prior assumptions. We briefly discuss some ways to do this below.

### 5.4.1 Uninformative priors

If we don't have strong beliefs about what  $\theta$  should be, it is common to use an **uninformative** or **non-informative** prior, and to let the data speak for itself.

### 5.4.2 Robust priors

In many cases, we are not very confident in our prior, so we want to make sure it does not have an undue influence on the result. This can be done by using **robust priors** (Insua and Ruggeri 2000), which typically have heavy tails, which avoids forcing things to be too close to the prior mean.

### 5.4.3 Mixtures of conjugate priors

Robust priors are useful, but can be computationally expensive to use. Conjugate priors simplify the computation, but are often not robust, and not flexible enough to encode our prior knowledge. However, it turns out that a **mixture of conjugate priors** is also conjugate, and can approximate any kind of prior (Dallal and Hall 1983; Diaconis and Ylvisaker 1985). Thus such priors provide a good compromise between computational convenience and flexibility.

## 5.5 Hierarchical Bayes

A key requirement for computing the posterior  $p(\theta|\mathcal{D})$  is the specification of a prior  $p(\theta|\eta)$ , where  $\eta$  are the hyper-parameters. What if we don't know how to set  $\eta$ ? In some cases, we can use uninformative priors, as we discussed above. A more Bayesian approach is to put a prior on our priors! In terms of graphical models (Chapter TODO), we can represent the situation as follows:

$$\eta \rightarrow \theta \rightarrow \mathcal{D} \quad (5.27)$$

This is an example of a **hierarchical Bayesian model**, also called a **multi-level model**, since there are multiple levels of unknown quantities.

## 5.6 Empirical Bayes

Method	Definition
Maximum likelihood	$\hat{\theta} = \arg \max_{\theta} p(\mathcal{D} \theta)$
MAP estimation	$\hat{\theta} = \arg \max_{\theta} p(\mathcal{D} \theta)p(\theta \eta)$
ML-II ( <b>Empirical Bayes</b> )	$\hat{\eta} = \arg \max_{\eta} \int p(\mathcal{D} \theta)p(\theta \eta)d\theta$ $= \arg \max_{\eta} p(\mathcal{D} \eta)$
MAP-II	$\hat{\eta} = \arg \max_{\eta} \int p(\mathcal{D} \theta)p(\theta \eta)p(\eta)d\theta$ $= \arg \max_{\eta} p(\mathcal{D} \eta)p(\eta)$
Full Bayes	$p(\theta, \eta \mathcal{D}) \propto p(\mathcal{D} \theta)p(\theta \eta)$

## 5.7 Bayesian decision theory

We have seen how probability theory can be used to represent and update our beliefs about the state of the world. However, ultimately our goal is to convert our beliefs into actions. In this section, we discuss the optimal way to do this.

Our goal is to devise a **decision procedure** or **policy**,  $f(x) : \mathcal{X} \rightarrow \mathcal{Y}$ , which minimizes the **expected loss**  $R_{\text{exp}}(f)$  (see Equation 1.2).

In the Bayesian approach to decision theory, the optimal output, having observed  $x$ , is defined as the output  $a$  that minimizes the **posterior expected loss**:

$$\rho(f) = \mathbb{E}_{p(y|x)}[L(y, f(x))] = \begin{cases} \sum_y L[y, f(x)]p(y|x) \\ \int_y L[y, f(x)]p(y|x)dy \end{cases} \quad (5.28)$$

Hence the **Bayes estimator**, also called the **Bayes decision rule**, is given by

$$\delta(x) = \arg \min_{f \in \mathcal{H}} \rho(f) \quad (5.29)$$

## 5.7.1 Bayes estimators for common loss functions

### 5.7.1.1 MAP estimate minimizes 0-1 loss

When  $L(y, f(x))$  is **0-1 loss**(Section 1.4.2.1), we can proof that MAP estimate minimizes 0-1 loss,

$$\begin{aligned}
 \arg \min_{f \in \mathcal{H}} \rho(f) &= \arg \min_{f \in \mathcal{H}} \sum_{i=1}^K L[C_k, f(\mathbf{x})] p(C_k | \mathbf{x}) \\
 &= \arg \min_{f \in \mathcal{H}} \sum_{i=1}^K \mathbb{I}(f(\mathbf{x}) \neq C_k) p(C_k | \mathbf{x}) \\
 &= \arg \min_{f \in \mathcal{H}} \sum_{i=1}^K p(f(\mathbf{x}) \neq C_k | \mathbf{x}) \\
 &= \arg \min_{f \in \mathcal{H}} [1 - p(f(\mathbf{x}) = C_k | \mathbf{x})] \\
 &= \arg \max_{f \in \mathcal{H}} p(f(\mathbf{x}) = C_k | \mathbf{x})
 \end{aligned}$$

### 5.7.1.2 Posterior mean minimizes $\ell_2$ (quadratic) loss

For continuous parameters, a more appropriate loss function is **squared error**,  $\ell_2$  **loss**, or **quadratic loss**, defined as  $L(y, f(\mathbf{x})) = [y - f(\mathbf{x})]^2$ .

The posterior expected loss is given by

$$\begin{aligned}
 \rho(f) &= \int_y L[y, f(\mathbf{x})] p(y | \mathbf{x}) dy \\
 &= \int_y [y - f(\mathbf{x})]^2 p(y | \mathbf{x}) dy \\
 &= \int_y [y^2 - 2yf(\mathbf{x}) + f(\mathbf{x})^2] p(y | \mathbf{x}) dy
 \end{aligned} \tag{5.30}$$

Hence the optimal estimate is the posterior mean:

$$\begin{aligned}
 \frac{\partial \rho}{\partial f} &= \int_y [-2y + 2f(\mathbf{x})] p(y | \mathbf{x}) dy = 0 \Rightarrow \\
 \int_y f(\mathbf{x}) p(y | \mathbf{x}) dy &= \int_y y p(y | \mathbf{x}) dy \\
 f(\mathbf{x}) \int_y p(y | \mathbf{x}) dy &= \mathbb{E}_{p(y | \mathbf{x})}[y] \\
 f(\mathbf{x}) &= \mathbb{E}_{p(y | \mathbf{x})}[y]
 \end{aligned} \tag{5.31}$$

This is often called the **minimum mean squared error** estimate or **MMSE** estimate.

### 5.7.1.3 Posterior median minimizes $\ell_1$ (absolute) loss

The  $\ell_2$  loss penalizes deviations from the truth quadratically, and thus is sensitive to outliers. A more robust alternative is the absolute or  $\ell_1$  loss. The optimal estimate is the posterior median, i.e., a value  $a$  such that  $P(y < a | \mathbf{x}) = P(y \geq a | \mathbf{x}) = 0.5$ .



*Proof.*

$$\begin{aligned}
 \rho(f) &= \int_y L[y, f(\mathbf{x})] p(y|\mathbf{x}) dy = \int_y |y - f(\mathbf{x})| p(y|\mathbf{x}) dy \\
 &= \int_y [f(\mathbf{x}) - y] p(y < f(\mathbf{x})|\mathbf{x}) + \\
 &\quad [y - f(\mathbf{x})] p(y \geq f(\mathbf{x})|\mathbf{x}) dy \\
 \frac{\partial \rho}{\partial f} &= \int_y [p(y < f(\mathbf{x})|\mathbf{x}) - p(y \geq f(\mathbf{x})|\mathbf{x})] dy = 0 \Rightarrow \\
 p(y < f(\mathbf{x})|\mathbf{x}) &= p(y \geq f(\mathbf{x})|\mathbf{x}) = 0.5 \\
 \therefore f(\mathbf{x}) &= \text{median}
 \end{aligned}$$

#### 5.7.1.4 Reject option

In classification problems where  $p(y|\mathbf{x})$  is very uncertain, we may prefer to choose a reject action, in which we refuse to classify the example as any of the specified classes, and instead say dont know. Such ambiguous cases can be handled by e.g., a human expert. This is useful in **risk averse** domains such as medicine and finance.

We can formalize the reject option as follows. Let choosing  $f(\mathbf{x}) = c_{K+1}$  correspond to picking the reject action, and choosing  $f(\mathbf{x}) \in \{C_1, \dots, C_k\}$  correspond to picking one of the classes. Suppose we define the loss function as

$$L(f(\mathbf{x}), y) = \begin{cases} 0 & \text{if } f(\mathbf{x}) = y \text{ and } f(\mathbf{x}), y \in \{C_1, \dots, C_k\} \\ \lambda_s & \text{if } f(\mathbf{x}) \neq y \text{ and } f(\mathbf{x}), y \in \{C_1, \dots, C_k\} \\ \lambda_r & \text{if } f(\mathbf{x}) = C_{K+1} \end{cases} \quad (5.32)$$

where  $\lambda_s$  is the cost of a substitution error, and  $\lambda_r$  is the cost of the reject action.

#### 5.7.1.5 Supervised learning

We can define the loss incurred by  $f(\mathbf{x})$  (i.e., using this predictor) when the unknown state of nature is  $\theta$  (the parameters of the data generating mechanism) as follows:

$$L(\theta, f) \triangleq \mathbb{E}_{p(\mathbf{x}, y|\theta)} [\ell(y - f(\mathbf{x}))] \quad (5.33)$$

This is known as the **generalization error**. Our goal is to minimize the posterior expected loss, given by

$$\rho(f|\mathcal{D}) = \int p(\theta|\mathcal{D}) L(\theta, f) d\theta \quad (5.34)$$

This should be contrasted with the frequentist risk which is defined in Equation TODO.

### 5.7.2 The false positive vs false negative tradeoff

In this section, we focus on binary decision problems, such as hypothesis testing, two-class classification, object/ event detection, etc. There are two types of error we can make: a **false positive** (aka **false alarm**), or a **false negative** (aka **missed detection**). The 0-1 loss treats these two kinds of errors equivalently. However, we can consider the following more general loss matrix:

TODO



# Chapter 6

## Linear Regression

### 6.1 Introduction

Linear regression is the work horse of statistics and (supervised) machine learning. When augmented with kernels or other forms of basis function expansion, it can model also nonlinear relationships. And when the Gaussian output is replaced with a Bernoulli or multinoulli distribution, it can be used for classification, as we will see below. So it pays to study this model in detail.

In the simplest approach, we can directly construct an appropriate function  $y(\mathbf{x})$  whose values for new inputs  $\mathbf{x}$  constitute the predictions for the corresponding values of  $t$ . More generally, from a probabilistic perspective, we aim to model the predictive distribution  $p(t|\mathbf{x})$  because this expresses the uncertainty about the value of  $t$  for each value of  $\mathbf{x}$ .

### 6.2 Representation

The simplest linear model for regression is **linear regression** that involves a linear combination of the input variables

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D \quad (6.1)$$

where  $\mathbf{x} = (x_1, \dots, x_D)^T$ .

Linear regression can be made to model non-linear relationships by replacing  $\mathbf{x}$  with some non-linear function of the inputs,  $\phi(\mathbf{x})$ . Consider the linear combinations of fixed nonlinear functions of the input variables, of the form

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x) = \sum_{j=0}^{M-1} w_j \phi_j(x) = \mathbf{w}^T \boldsymbol{\phi}(x) \quad (6.2)$$

where  $\phi_j(x)$  are known as **basis functions**,  $\mathbf{w} = (w_0, \dots, w_{M-1})^T$  and  $\boldsymbol{\phi} = (\phi_0, \dots, \phi_{M-1})^T$ . This is known as **basis function expansion**. (Note that the model is still linear in the parameters  $\mathbf{w}$ , so it is still called linear regression; the importance of this will become clear below.) A simple example are polynomial basis functions, where the model has the form

$$\boldsymbol{\phi}(x) = (1, x, \dots, x^d) \quad (6.3)$$

As before, we assume that the target variable is given by a deterministic function  $y(\mathbf{x}, \mathbf{w})$  with additive Gaussian noise so that

$$t = y(\mathbf{x}, \mathbf{w}) + \varepsilon \quad (6.4)$$

where  $\varepsilon$  is a zero mean Gaussian random variable with precision (inverse variance)  $\beta$ . Thus

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|h(\mathbf{x}, \mathbf{w}), \beta^{-1}) \quad (6.5)$$

or

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \varepsilon \quad (6.6)$$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mathbf{w}^T \mathbf{x}, \sigma^2) \quad (6.7)$$

$$(6.8)$$

where  $\mathbf{w}$  (**weight vector**) and  $\mathbf{x}$  are extended vectors,  $\mathbf{x} = (1, x)$ ,  $\mathbf{w} = (b, w)$  and  $\varepsilon$  has a **Gaussian** or **normal** distribution.  $w_0$  is the intercept or **text**

### 6.3 Maximum likelihood estimations(least squares)

Assume the training examples are **independently and identically distributed(IID)**,we obtain the **likelihood function**,which is a function of adjustable parameters  $w$  and  $\beta$ ,in the form

$$p(t|X, w, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | w^T \phi(x_n), \beta^{-1}) \quad (6.9)$$

we can write the **log-likelihood** (logarithm of the likelihood) function as follows:

$$\ell(\theta) \triangleq \log p(\mathcal{D}|\theta) \quad (6.10)$$

A common way to estimate the parameters of a statistical model is to compute the MLE(Maximum likelihood estimations),which is defined as

$$\hat{\theta} = \arg \max_{\theta} \log p(\mathcal{D}|\theta) \quad (6.11)$$

Instead of maximizing the log-likelihood, we can equivalently minimize the **negative log likelihood** or **NLL**:

$$\text{NLL}(\theta) \triangleq -\ell(\theta) = -\log p(\mathcal{D}|\theta) \quad (6.12)$$

The NLL formulation is sometimes more convenient, since many optimization software packages are designed to find the minima of functions, rather than maxima.

Now let us apply the method of MLE to the linear regression setting. Inserting the definition of the Gaussian into the above, we find that the log likelihood is given by

$$\ell(\theta) = \log p(T|w, \beta) \quad (6.13)$$

$$= \sum_{n=1}^N \log \mathcal{N}(t_n | w^T \phi(x_n), \beta^{-1}) \quad (6.14)$$

$$= \sum_{i=1}^N \log \left[ \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right) \right] \quad (6.15)$$

$$= \frac{N}{2} \log \beta - \frac{N}{2} \log(2\pi) - \beta E_D(w) \quad (6.16)$$

$$= -\frac{1}{2\sigma^2} \text{RSS}(w) - \frac{w}{2} \log(2\pi\sigma^2) \quad (6.17)$$

where the sum-of-squares error function is defined by

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2 \quad (6.18)$$

RSS stands for **residual sum of squares** and is defined by

$$\text{RSS}(w) \triangleq \sum_{i=1}^N (t_i - w^T x_i)^2 \quad (6.19)$$

The RSS can also be written as the square of the  $\ell_2$  **norm** of the vector of residual errors:

$$\text{RSS}(w) = \|\varepsilon\|_2^2 = \sum_{i=1}^N \varepsilon_i^2 \quad (6.20)$$

where  $\varepsilon_i = (t_i - w^T x_i)$ .

There two ways to compute the optimal solution of likelihood function.

### 6.3.1 Derivations of the MLE

We see that maximizing the likelihood function under a conditional Gaussian noise distribution for a linear model is equivalent to  $\mathbf{w}$  minimizing the RSS(or  $E_D$ ), so this method is known as **least squares**. The gradient of the log likelihood function takes the form

$$\nabla \log p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T \quad (6.21)$$

Setting the gradient respect to  $\mathbf{w}$  to zero gives

$$0 = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left( \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right) \quad (6.22)$$

Solving for  $\mathbf{w}$  we obtain

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (6.23)$$

which are known as the **normal equations**.

Make the bias parameter explicit, the error function

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n)\}^2 \quad (6.24)$$

Setting the derivative with respect to  $w_0$  equal to zero, we obtain

$$w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j \quad (6.25)$$

Then we have

$$y = \mathbf{w}^T \phi \quad (6.26)$$

$$y = \sum_{j=1}^M w_j \phi_j + w_0 \quad (6.27)$$

$$y = \sum_{j=1}^M w_j \phi_j + \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j \quad (6.28)$$

$$y - \bar{y} = \sum_{j=1}^M w_j (\phi_j - \bar{\phi}_j) \quad (6.29)$$

which indicates that we can normalize the features by subtracting the mean of  $\Phi$  while training to reduce computation complexity but obtain the same weight parameter.

Note that  $\mathbf{t} = (t_1, t_2, \dots, t_N)^T$ ,  $\Phi = \begin{pmatrix} \phi_1^T \\ \phi_2^T \\ \vdots \\ \phi_N^T \end{pmatrix}$

$\Phi$  is an  $N \times M$  matrix, called the **design matrix**, whose elements are given by  $\Phi_{nj} = \phi_j(\mathbf{x}_n)$  so that

$$\Phi = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \dots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{bmatrix} \quad (6.30)$$

*Proof.*

$$\therefore \sum_{n=1}^N t_n \phi(x_n)^T = \mathbf{t}^T \Phi \quad (6.31)$$

$$\left\{ \begin{array}{l} \Phi = \begin{bmatrix} \Phi_1^T \\ \Phi_2^T \\ \dots \\ \Phi_n^T \end{bmatrix} \\ \Phi^T \Phi = \begin{bmatrix} \Phi_1 & \Phi_2 & \dots & \Phi_n \end{bmatrix} \cdot \begin{bmatrix} \Phi_1^T \\ \Phi_2^T \\ \dots \\ \Phi_n^T \end{bmatrix} = \sum_{i=1}^N \phi_i \phi_i^T = \text{sum of squares matrix} \end{array} \right. \quad (6.32)$$

$$\therefore \mathbf{w}^T = \mathbf{t}^T \Phi (\Phi^T \Phi)^{-1} \therefore \mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (6.33)$$

*Proof.* Let's drop constants with respect to  $\mathbf{w}$  and NLL can be written as

$$\text{NLL}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^T \phi_i)^2 \quad (6.34)$$

Then we can rewritten the objective in a form that is more amenable to differentiation:

$$\text{NLL}(\mathbf{w}) = \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) \quad (6.35)$$

$$= \frac{1}{2} [\mathbf{w}^T (\Phi^T \Phi) \mathbf{w} - \mathbf{t}^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T \mathbf{t} + \mathbf{t}^T \mathbf{t}] \quad (6.36)$$

$$\Rightarrow \frac{\partial}{\partial \mathbf{w}} \text{NLL} = \frac{1}{2} \mathbf{w}^T (\Phi^T \Phi) \mathbf{w} - \mathbf{w}^T (\Phi^T \mathbf{t}) \quad (6.37)$$

where

When  $\mathcal{D}$  is small (for example,  $N < 1000$ ), we can use the following equation to compute  $\mathbf{w}$  directly

$$\mathbf{w}_{ML} = \hat{\mathbf{w}}_{\text{OLS}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (6.38)$$

The corresponding solution  $\hat{\mathbf{w}}_{\text{OLS}}$  to this linear system of equations is called the **ordinary least squares** or **OLS** or **normal equation** solution.

*Proof.* We now state without proof some facts of matrix derivatives (we wont need all of these at this section).

$$\text{tr} A \triangleq \sum_{i=1}^n A_{ii}$$

$$\frac{\partial}{\partial A} AB = B^T \quad (6.39)$$

$$\frac{\partial}{\partial A^T} f(A) = \left[ \frac{\partial}{\partial A} f(A) \right]^T \quad (6.40)$$

$$\frac{\partial}{\partial A} ABA^T C = CAB + C^T AB^T \quad (6.41)$$

$$\frac{\partial}{\partial A} |A| = |A| (A^{-1})^T \quad (6.42)$$

Then,

$$\begin{aligned}
\text{NLL}(\mathbf{w}) &= \frac{1}{2N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \\
\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y}) \\
&= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w}) \\
&= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} \\
&= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} (\text{tr} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2 \text{tr} \mathbf{y}^T \mathbf{X} \mathbf{w})
\end{aligned}$$

Combining Equations 6.40 and 6.41, we find that

$$\frac{\partial}{\partial \mathbf{A}^T} \mathbf{A} \mathbf{B} \mathbf{A}^T \mathbf{C} = \mathbf{B}^T \mathbf{A}^T \mathbf{C}^T + \mathbf{B} \mathbf{A}^T \mathbf{C}$$

Let  $\mathbf{A}^T = \mathbf{w}$ ,  $\mathbf{B} = \mathbf{B}^T = \mathbf{X}^T \mathbf{X}$ , and  $\mathbf{C} = \mathbf{I}$ , Hence,

$$\begin{aligned}
\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= \frac{1}{2} (\mathbf{X}^T \mathbf{X} \mathbf{w} + \mathbf{X}^T \mathbf{X} \mathbf{w} - 2 \mathbf{X}^T \mathbf{y}) \\
&= \frac{1}{2} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) \\
\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= 0 \Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = 0 \\
\mathbf{X}^T \mathbf{X} \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\
\hat{\mathbf{w}}_{\text{OLS}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}
\end{aligned} \tag{6.43}$$

Equation 6.43 is known as the **normal equation**.

*Proof.* Consider **vector – vector** differentiation, if  $\mathbf{A}$  is not a function of  $\mathbf{x}$  then

$$\frac{\partial \mathbf{X}^T \mathbf{A} \mathbf{X}}{\partial \mathbf{X}} = \mathbf{X}^T (\mathbf{A} + \mathbf{A}^T) = (\mathbf{A} + \mathbf{A}^T) \mathbf{X} \tag{6.44}$$

Then, we can derive  $\frac{\partial \text{NLL}}{\partial \mathbf{w}}$

$$g(\mathbf{w}) = \frac{\partial \text{NLL}}{\partial \mathbf{w}} \tag{6.45}$$

$$= \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{y} \tag{6.46}$$

Setting the derivative to zero, we get

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \tag{6.47}$$

### 6.3.2 Geometric interpretation

See Figure 6.1. Given that

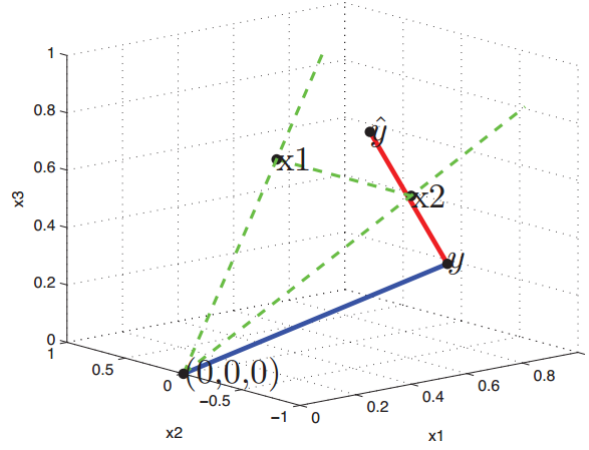


Fig. 6.1: Graphical interpretation of least squares for  $N = 3$  examples and  $D = 2$  features.  $\tilde{x}_1$  and  $\tilde{x}_2$  are vectors in  $\mathbb{R}^3$ ; together they define a 2D plane.  $y$  is also a vector in  $\mathbb{R}^3$  but does not lie on this 2D plane. The orthogonal projection of  $y$  onto this plane is denoted  $\hat{y}$ . The red line from  $y$  to  $\hat{y}$  is the residual, whose norm we want to minimize. For visual clarity, all vectors have been converted to unit norm.

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \dots \\ x_N^T \end{bmatrix} = [\tilde{x}_1 \ \tilde{x}_1 \ \dots \ \tilde{x}_D] \quad (6.48)$$

$$y = \begin{bmatrix} t_1 \\ t_2 \\ \dots \\ t_n \end{bmatrix} \quad (6.49)$$

We seek a vector  $\hat{t} \in \mathbb{R}^N$  that lies in the column linear space of  $X$  and is as close as possible to  $t$ , i.e., we want to find

$$\hat{t} \in \text{span}(X) \quad (6.50)$$

$$\Rightarrow \hat{t} = Xw = w_1\tilde{x}_1 + \dots + w_D\tilde{x}_D \quad (6.51)$$

$$\hat{t} = \arg \min_{\hat{t} \in \text{span}(\{\tilde{x}_1, \dots, \tilde{x}_D\})} \quad (6.52)$$

To minimize the norm of the residual,  $t - \hat{t}$ , we want the residual vector to be orthogonal to every column of  $X$ , so  $\tilde{x}_j(t - \hat{t}) = 0$  for  $j = 1 : D$ . Hence

$$\begin{aligned} \tilde{x}_j(t - \hat{t}) = 0 &\Rightarrow X^T(t - Xw) = 0 \\ &\Rightarrow w = (X^T X)^{-1} X^T t \end{aligned} \quad (6.53)$$

### 6.3.3 Sequential learning

Batch techniques, such as the maximum likelihood solution, which involves processing the entire training set in one go (pass), can be computationally costly for large data sets. If the data set is sufficiently large, it may be worthwhile to use **sequential** algorithms, known as **on-line** algorithms. When  $D$  is large, use **stochastic gradient descent (SGD)**, also known as sequential gradient descent.



$$E = \sum_n E_n \quad (6.54)$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_n \quad (6.55)$$

where  $\tau$  denotes the iteration number, the  $\eta$  is a learning rate parameter. For the case of sum-of-squares error function, this gives

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta (t_n - \mathbf{w}^{(\tau)T} \phi_n) \phi_n \quad (6.56)$$

where  $\phi_n = \phi(\mathbf{x}_n)$ . This is known as the **least-mean-squares** or the LMS algorithm.

$$\therefore \frac{\partial}{\partial w_i} \text{NLL}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i) x_{ij} \quad (6.57)$$

$$\begin{aligned} \therefore w_j &= w_j - \alpha \frac{\partial}{\partial w_j} \text{NLL}(\mathbf{w}) \\ &= w_j - \sum_{i=1}^N \alpha (\mathbf{w}^T \mathbf{x}_i - y_i) x_{ij} \end{aligned} \quad (6.58)$$

$$\therefore \mathbf{w} = \mathbf{w} - \alpha (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x} \quad (6.59)$$

## 6.4 Ridge regression(MAP)

One problem with ML estimation is that it can result in over-fitting. In this section, we discuss a way to ameliorate this problem by using MAP estimation with a Gaussian prior.

### 6.4.1 Basic idea

We can encourage the parameters to be small, thus resulting in a smoother curve, by using a zero-mean Gaussian prior:

$$p(\mathbf{w}) = \prod_j \mathcal{N}(w_j | 0, \tau^2) \quad (6.60)$$

where  $1/\tau^2$  controls the strength of the prior. The corresponding MAP estimation problem becomes

$$\arg \max_{\mathbf{w}} \sum_{i=1}^N \log \mathcal{N}(t_i | w_0 + \mathbf{w}^T \phi_i, \sigma^2) + \sum_{j=1}^D \log \mathcal{N}(w_j | 0, \tau^2) \quad (6.61)$$

This is equivalent to minimizing the following

$$J(\mathbf{w}) = E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) \quad (6.62)$$

$$= \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \mathbf{w}^T \mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|^2, \lambda \triangleq \frac{\sigma^2}{\tau^2} \quad (6.63)$$

where the first term is the MSE/NLL as usual, and the second term,  $\lambda \geq 0$ , is the regularization coefficient that controls the complexity penalty. Here  $E_W(\mathbf{w})$  is one of the simplest forms of regularizer given by the sum-of-squares of the weight vector elements.

$$E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad (6.64)$$

This particular choice of regularizer is known as **weight decay**.

The corresponding solution is given by

$$\hat{\mathbf{w}}_{\text{ridge}} = (\lambda \mathbf{I}_D + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (6.65)$$

This technique is known as **ridge regression**, or **penalized least squares**. In general, adding a Gaussian prior to the parameters of a model to encourage them to be small is called  $\ell_2$  **regularization** or **weight decay**. Note that the offset term  $w_0$  is not regularized, since this just affects the height of the function, not its complexity.

We will consider a variety of different priors in this book. Each of these corresponds to a different form of **regularization**. This technique is very widely used to prevent overfitting.

## 6.4.2 Multiple outputs

Now consider the case where we wish to predict  $K > 1$  target variables, which we denote collectively by the target vector  $\mathbf{t}$ . To use the same set of basis functions to model all the components of the target vector so that

$$\mathbf{y}(x, \mathbf{w}) = \mathbf{W}^T \phi(x) \quad (6.66)$$

where  $\mathbf{y}$  is a  $K$ -dimensional column vector,  $\mathbf{W}$  is an  $M \times K$  matrix of parameters.

The conditional distribution of the target vector is an isotropic Gaussian

$$p(\mathbf{t} | \mathbf{x}, \mathbf{W}, \beta) = \mathcal{N}(\mathbf{t} | \mathbf{W}^T \phi(x), \beta^{-1} \mathbf{I}) \quad (6.67)$$

The log likelihood function is then given by

$$\ln p(\mathbf{T} | \mathbf{X}, \mathbf{W}, \beta) = \sum_{n=1}^N \ln \mathcal{N}(\mathbf{t}_n | \mathbf{W}^T \phi(x_n), \beta^{-1} \mathbf{I}) \quad (6.68)$$

$$= \frac{NK}{2} \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} \sum_{n=1}^N \|\mathbf{t}_n - \mathbf{W}^T \phi(x_n)\|^2 \quad (6.69)$$

We maximize this function with respect to  $\mathbf{W}$ , giving

$$\mathbf{W}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T} \quad (6.70)$$

## 6.4.3 Numerically stable computation \*

$$\hat{\mathbf{w}}_{\text{ridge}} = \mathbf{V}(\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^T \mathbf{y} \quad (6.71)$$

## 6.4.4 Connection with PCA \*

## 6.4.5 Regularization effects of big data

Regularization is the most common way to avoid overfitting. However, another effective approach which is not always available is to use lots of data. It should be intuitively obvious that the more training data we have, the better we will be able to learn.

In domains with lots of data, simple methods can work surprisingly well (Halevy et al. 2009). However, there are still reasons to study more sophisticated learning methods, because there will always be problems for which we have little data. For example, even in such a data-rich domain as web search, as soon as we want to start personalizing the results, the amount of data available for any given user starts to look small again (relative to the complexity of the problem).

## 6.5 The Bias-Variance Decomposition

### 6.5.0.1 representation

Data representation:  $N$  observations of  $\mathbf{x}$ , written  $X \equiv (x_1, x_2, \dots, x_N)^T$  together with corresponding observations of the values of  $t$ , denoted  $t \equiv (t_1, t_2, \dots, t_N)^T$ .

$$t_i = f(\mathbf{x}_i) + \varepsilon \quad (6.72)$$

where the noise  $\varepsilon$  has zero mean and variance  $\sigma^2$ . Find a function

$$\hat{f}(x)$$

that approximates the true function

$$t = f(\mathbf{x})$$

as well as possible. Make "as well as possible" precise by measuring the mean squared error between  $y$  and  $\hat{f}(x)$ , we want  $(t - \hat{f}(x))^2$  to be minimal. Hypothesis function (model representation):  $y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$ .  $M$  is the order of the polynomial, and  $x^j$  denotes  $x$  raised to the power of  $j$ . The polynomial coefficients  $w_0, \dots, w_M$  are collectively denoted by the vector  $\mathbf{w}$ .

Error Function (Sum of squares of errors between predictions  $y(x_n, \mathbf{w})$  for each data point  $x_n$  and the corresponding target values  $t_n$ , so that we minimize:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \quad (6.73)$$

root-mean-square (RMS) error defined by

$$E_{RMS} = \sqrt{2E(\mathbf{w}^*)/N} \quad (6.74)$$

Penalized (regularized) error function

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{1}{2} \|\mathbf{w}\|^2 \quad (6.75)$$

where  $\|\mathbf{w}\|^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + \dots + w_M^2$

### 6.5.0.2 Loss function for regression

$$\mathbb{E}[L] = \iint L(t, y(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt \quad (6.76)$$

A common choice of loss function in squared loss given by

$$L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2 \quad (6.77)$$

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt. \quad (6.78)$$

Minimize  $\mathbb{E}[L]$  by using the calculus of variations to give

$$\frac{\delta \mathbb{E}[L]}{\delta y(\mathbf{x})} = 2 \int \{y(\mathbf{x}) - t\} p(\mathbf{x}, t) dt = 0 \quad (6.79)$$

Solving for  $y(\mathbf{x})$  and using the sum and product rules of probability, we obtain

$$y(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int t p(t|\mathbf{x}) dt = \mathbb{E}[t|\mathbf{x}] \quad (6.80)$$

Let's derive this result in a slightly different way. Armed with knowledge that the optimal solution is the conditional expectation, we can expand the square term as follows

$$\{y(\mathbf{x}-t)\}^2 = \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 = \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2 \quad (6.81)$$

where,  $\mathbb{E}[t|\mathbf{x}]$  denote  $\mathbb{E}_t[t|\mathbf{x}]$ . Substitute into the loss function and perform the integral over  $t$ , we see the cross-term vanishes

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \quad (6.82)$$

$$= \int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{\mathbb{E}[t|\mathbf{x}] - t\}^2 p(\mathbf{x}) d\mathbf{x} \quad (6.83)$$

$$= \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}) d\mathbf{x} \quad (6.84)$$

$$= \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \quad (6.85)$$

$$(6.86)$$

### 6.5.0.3 Decomposition

For a popular choice, we use squared loss function, for which the optimal prediction is given by the conditional expectation, which we denote by  $h(\mathbf{x})$  and which is given by

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt \quad (6.87)$$

Consider the integrand of the first term of 6.82, which for particular data set  $D$  takes the form

$$\{y(\mathbf{x}; D) - h(\mathbf{x})\}^2 \quad (6.88)$$

This quantity will be dependent on the particular data set  $D$ , so we take its average over the ensemble of data sets. If we add and subtract the quantity  $\mathbb{E}_D[y(\mathbf{x}; D)]$  inside the braces, and then expand, we obtain

$$\{y(\mathbf{x}; D) - h(\mathbf{x})\}^2 \quad (6.89)$$

$$= \{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)] + \mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2 \quad (6.90)$$

$$= \{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)]\}^2 + \{\mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2 + 2\{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)]\}\{\mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\} \quad (6.91)$$

$$= \underbrace{\{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)]\}^2}_{\text{variance}} + \underbrace{\{\mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + 0 \quad (6.92)$$

$$(6.93)$$

The decomposition of the expected squared loss

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise} \quad (6.94)$$

where

$$(\text{bias})^2 = \dots \quad (6.95)$$

$$\text{variance} = \dots \quad (6.96)$$

$$\text{noise} = \dots \quad (6.97)$$

The function  $y(\mathbf{x})$  we seek to determine enters only the first term, which will be minimized when  $y(\mathbf{x})$  is equal to  $\mathbb{E}[t|\mathbf{x}]$ , in which case this term will vanish. The second term is the variance of distribution of  $t$ , averaged over  $\mathbf{x}$ , representing the intrinsic variability of the target data and can be regarded as noise. It's the irreducible minimum value of the loss function.

More sophisticated loss function, Minkowski loss

$$\mathbb{E}[L_q] = \iint |y(\mathbf{x}) - t|^q p(\mathbf{x}, t) d\mathbf{x} dt \quad (6.98)$$

## 6.6 Bayesian linear regression

Hold-out data can be used to determine model complexity but it will be computationally expensive and wasteful of valuable data. We therefore turn to a Bayesian treatment of linear regression, which will avoid the over-fitting problem of maximum likelihood, and which will also lead to automatic methods of determining model complexity using the training data.

### 6.6.1 Parameter distribution

First introduce a prior probability distribution over the model parameters. The likelihood function  $p(t|\mathbf{w})$  defined by 6.9 is the exponential of a quadratic function of  $\mathbf{w}$ , so the corresponding conjugate prior is Gaussian

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0) \quad (6.99)$$

The posterior distribution is proportional to the product of the likelihood and the prior. And the posterior will also be Gaussian due to the choice of conjugate Gaussian prior, which is derived in 2.5.

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \quad (6.100)$$

where

$$\mathbf{m}_N = \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\Phi^T\mathbf{t}) \quad (6.101)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta\Phi^T\Phi \quad (6.102)$$

Thus the maximum posterior weight vector is simply given by  $\mathbf{w}_{MAP} = \mathbf{m}_N$ . If  $N = 0$  then the posterior distribution reverts to the prior. Furthermore, if data points arrive sequentially, then the posterior distribution at any stage acts as the prior distribution, such that the new posterior is again given.

A zero-mean isotropic Gaussian governed by a single precision parameter  $\alpha$  so that

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \quad (6.103)$$

for which the posterior is given by

$$\mathbf{m}_N = \beta\mathbf{S}_N\Phi^T\mathbf{t} \quad (6.104)$$

$$\mathbf{S}_N^{-1} = \alpha\mathbf{I} + \beta\Phi^T\Phi \quad (6.105)$$

The log of posterior distribution is given by the sum of the log likelihood and the log of prior and, as a function of  $\mathbf{w}$ , takes the form

$$\log p(\mathbf{w}|\mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(x_n)\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const} \quad (6.106)$$

Maximization of this posterior w.r.t  $\mathbf{w}$  is equivalent to minimization of the sum-of-squares error function with the addition of a quadratic regularization term, corresponding with  $\lambda = \alpha/\beta$ .

### 6.6.2 Predictive distribution

Evaluate the **predictive distribution** defined by

$$p(t|\mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \quad (6.107)$$

The conditional distribution  $p(t|\mathbf{x}, \mathbf{w}, \beta)$  of the target variable is given by 6.5, and the posterior weight distribution is given by 6.100. This involves the **convolution** of two Gaussian distributions. The predictive distribution takes the form

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \quad (6.108)$$

where the variance  $\sigma_N^2(\mathbf{x})$  is given by

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}) \quad (6.109)$$

The first term represents the noise on the data whereas the second term reflects the uncertainty associated with the parameters  $\mathbf{w}$ .  $\sigma_{N+1}^2(\mathbf{x}) \leq \sigma_N^2(\mathbf{x})$ . In the limit  $N \rightarrow \infty$ , the second term goes to zero.

Note that if both  $\mathbf{w}$  and  $\beta$  are treated as unknown, then we can introduce a conjugate prior distribution  $p(\mathbf{w}, \beta)$  given by Gaussian-gamma distribution, leading to a Student's t-distribution predictive distribution.

### 6.6.3 Equivalent kernel

The posterior mean solution ?? has an interpretation that will set stage for kernel methods, including Gaussian processes. The predictive mean:

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{m}_N^T \phi(\mathbf{x}) = \beta \phi(\mathbf{x})^T \mathbf{S}_N \Phi^T \mathbf{t} = \sum_{n=1}^N \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_n) t_n = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n \quad (6.110)$$

where

$$k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}') \quad (6.111)$$

is known as the **smoother matrix** or **equivalent kernel**. Regression functions, such as this, which make predictions by taking linear combinations of the training set target values are known **linear smoothers**. The kernel functions are localized around  $\mathbf{x}$  (local evidence weight more than distant evidence).

The covariance between  $y(\mathbf{x})$  and  $y(\mathbf{x}')$  is given by

$$\text{cov}[y(\mathbf{x}), y(\mathbf{x}')] = \text{cov}[\phi(\mathbf{x})^T \mathbf{w}, \mathbf{w}^T \phi(\mathbf{x}')] \quad (6.112)$$

$$= \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}') = \beta^{-1} k(\mathbf{x}, \mathbf{x}') \quad (6.113)$$

We see that the predictive mean at nearby points will be highly correlated.

The formulation of linear regression in terms of kernel functions suggests an alternative approach, called **Gaussian process**: Instead of introducing a set of basis functions, which implicitly determines an equivalent kernel, we can instead define a localized kernel directly and use this to make predictions for new input vectors  $\mathbf{x}$ , given the observed training set.

The effective kernel defines the weights by which the training set target values are combined in order to make a prediction at a new value of  $\mathbf{x}$ , and it can be shown that these weights sum to one,

$$\sum_{n=1}^N \alpha_n = 1 \quad (6.114)$$

for all values of  $\mathbf{x}$ .

General kernel functions share an important property, namely that it can be expressed in the form an inner product with respect to a vector  $\Psi(\mathbf{x})$  of nonlinear functions, so that

$$k(\mathbf{x}, \mathbf{z}) = \Psi(\mathbf{x})^T \Psi(\mathbf{z}) \quad (6.115)$$

where  $\Psi(\mathbf{x}) = \beta^{1/2} \mathbf{S}_N^{1/2} \phi(\mathbf{x})$

## 6.7 Bayesian Model Comparison

The over-fitting associated with maximum likelihood can be avoided by **marginalizing(summing or integrating)** over the model parameters instead of making point estimates of their values, no need for validation.

Suppose we wish to compare a set of  $L$  models  $\{\mathcal{M}_i\}, i = 1, \dots, L$  over observed data  $\mathcal{D}$ . Evaluate the posterior distribution

$$p(\mathcal{M}_i|\mathcal{D}) \propto p(\mathcal{M}_i)p(\mathcal{D}|\mathcal{M}_i) \quad (6.116)$$

The prior allows us to express a preference for different models.  $p(\mathcal{D}|\mathcal{M}_i)$  is the **model evidence**, which expresses the preference shown by the data for different models. It is also called the **marginal likelihood**. The ratio of model evidences  $p(\mathcal{D}|\mathcal{M}_i)/p(\mathcal{D}|\mathcal{M}_j)$  for two models is **Bayes factor**. The predictive distribution

$$p(t|\mathbf{x}, \mathcal{D})p(\mathcal{M}_i|\mathcal{D}) = p(t|\mathbf{x}, \mathcal{M}_i, \mathcal{D})p(\mathcal{M}_i|\mathcal{D}) \quad (6.117)$$

$$\Rightarrow \sum_{i=1}^L p(t|\mathbf{x}, \mathcal{D})p(\mathcal{M}_i|\mathcal{D}) = \sum_{i=1}^L p(t|\mathbf{x}, \mathcal{M}_i, \mathcal{D})p(\mathcal{M}_i|\mathcal{D}) \quad (6.118)$$

$$\Rightarrow p(t|\mathbf{x}, \mathcal{D}) = \sum_{i=1}^L p(t|\mathbf{x}, \mathcal{M}_i, \mathcal{D})p(\mathcal{M}_i|\mathcal{D}) \quad (6.119)$$

$$(6.120)$$

This is an example of a **mixture distribution** in which the overall predictive distribution is obtained by averaging the predictive distributions of individual models weighted by the posterior probabilities  $p(\mathcal{M}_i|\mathcal{D})$  of those models.

To use the single most probable model alone to make predictions is known as **model selection**. Model evidence

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathcal{M}_i)d\mathbf{w} \quad (6.121)$$

because

$$p(\mathbf{w}|\mathcal{D}, \mathcal{M}_i) = \frac{p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)p(\mathbf{w}|\mathcal{M}_i)}{p(\mathcal{D}|\mathcal{M}_i)} \quad (6.122)$$

Assume the posterior distribution over parameters is sharply peaked around the most probable value  $w_{MAP}$ , with width  $\Delta w_{posterior}$ , then we can approximate the integral by the value of integrand at its maximum times the width of the peak. And further assume that the prior is flat with width  $\Delta w_{prior}$  so that  $p(w) = 1/\Delta w_{prior}$ , then we have

$$p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w} \simeq p(\mathcal{D}|\mathbf{w}_{MAP})p(\mathbf{w}_{MAP})\Delta w_{posterior} \simeq p(\mathcal{D}|\mathbf{w}_{MAP})\frac{\Delta w_{posterior}}{\Delta w_{prior}} \quad (6.123)$$

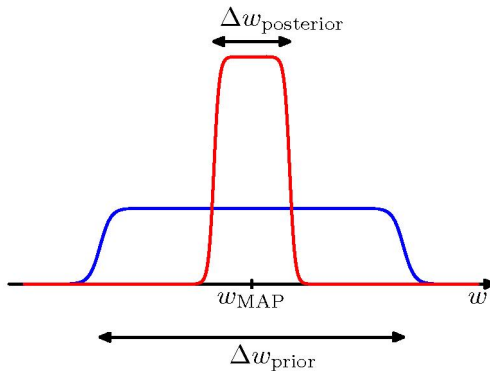


Fig. 6.2: model evidence approximation

and so taking logs

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|w_{MAP}) + \ln\left(\frac{\Delta w_{posterior}}{\Delta w_{prior}}\right) \quad (6.124)$$

If we have a set of  $M$  parameters, assuming that all parameters have the same ratio of  $\Delta w_{posterior}/\Delta w_{prior}$ , we obtain

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|w_{MAP}) + M \ln\left(\frac{\Delta w_{posterior}}{\Delta w_{prior}}\right) \quad (6.125)$$

As we increase the complexity of the model, the first term representing the fit to the data, increases(?), whereas the second term will decrease due to the dependence on  $M$ . The optimal model is given by a trade-off between these two competing terms. Average the Bayes factor over the distribution of data sets

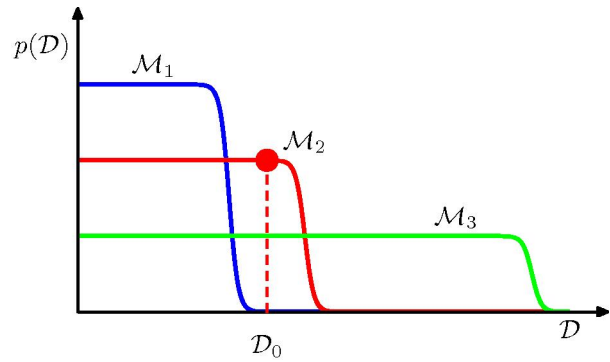


Fig. 6.3: Model complexity increases

$$\int p(\mathcal{D}|\mathcal{M}_1) \ln \frac{p(\mathcal{D}|\mathcal{M}_1)}{p(\mathcal{D}|\mathcal{M}_2)} \quad (6.126)$$

This quantity is *Kullback – Leibler* divergence, and will be bigger than or equal to 0 if  $\mathcal{M}_1$  is the correct model.



## Chapter 7

# Linear Models for Classification

The goal in classification is to take an input vector  $\mathbf{x}$  and to assign it to one of  $K$  discrete classes  $\mathcal{C}_k$  where  $K = 1, \dots, K$ . The input space is thereby divided into **decision regions** whose boundaries are called **decision boundaries** or **decision surfaces**. Data sets whose classes can be separated exactly by linear decision surfaces are said to be **linearly separable**.

In Chapter 1, we identified three distinct approaches to the classification problem. The simplest involves constructing a **discriminant function** that directly assigns each vector  $\mathbf{x}$  to a specific class. A more powerful approach, however, models the **conditional probability distribution**  $p(\mathcal{C}_k|\mathbf{x})$  in an **inference** stage, and then subsequently uses this distribution to make optimal **decisions**. By separating inference and decision, we gain numerous benefits, as discussed in 1. There are two different approaches to determining the conditional probabilities  $p(\mathcal{C}_k|\mathbf{x})$ . One technique is to model them **directly**, for example by representing them as parametric models and then optimizing the parameters using a training set. Alternatively, we can adopt a generative approach in which we model the **class-conditional densities** given by  $p(\mathbf{x}|\mathcal{C}_k)$ , together with the **prior probabilities**  $p(\mathcal{C}_k)$  for the classes, and then we compute the required posterior probabilities using Bayes theorem

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} \quad (7.1)$$

In the linear regression models considered in Chapter 3, the model prediction  $y(\mathbf{x}, \mathbf{w})$  was given by a linear function of the parameters  $\mathbf{w}$ . In the simplest case, the model is also linear in the input variables and therefore takes the form  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ , so that  $y$  is a real number. For classification problems, however, we wish to predict discrete class labels, or more generally posterior probabilities that lie in the range  $(0, 1)$ . To achieve this, we consider a generalization of this model in which we transform the linear function of  $\mathbf{w}$  using a nonlinear function  $f(\cdot)$  so that

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0) \quad (7.2)$$

In the machine learning literature  $f(\cdot)$  is known as an **activation function**, whereas its inverse is called a **link function** in the statistics literature. The decision surfaces correspond to  $y(\mathbf{x}) = \text{constant}$ , so that  $\mathbf{w}^T \mathbf{x} + w_0 = \text{constant}$  and hence the decision surfaces are linear functions of  $\mathbf{x}$ , even if the function  $f(\cdot)$  is nonlinear. For this reason, the class of models described by (4.3) are called **generalized linear models** (McCullagh and Nelder, 1996). However, in contrast to the models used for regression, they are no longer linear in the parameters due to the presence of the nonlinear function  $f(\cdot)$ .

## 7.1 Discriminant Functions

A discriminant is a function that takes an input vector  $\mathbf{x}$  and assigns it to be one of  $K$  classes, denoted  $\mathcal{C}_k$ .

### 7.1.1 Two classes

Linear discriminant function

$$y(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + w_0 \quad (7.3)$$

where  $\mathbf{w}$  is called a **weight vector**, and  $w_0$  is *bias*. The negative of the bias is sometimes called **threshold**.

The value of  $y(\mathbf{x})$  gives a signed measure of the perpendicular distance  $r$  of point  $\mathbf{x}$  from the decision surface. Consider an arbitrary point  $\mathbf{x}$  and let  $\mathbf{x}_\perp$  be its orthogonal projection onto the decision surface, then

$$\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \quad (7.4)$$

$$(7.5)$$

Multiplying both sides by  $\mathbf{w}^T$  and adding  $w_0$ ,

$$y(\mathbf{x}) = y(\mathbf{x}_\perp) + r \frac{\mathbf{w} \cdot \mathbf{w}}{\|\mathbf{w}\|} \quad (7.6)$$

$$= r \|\mathbf{w}\| \quad (7.7)$$

$$\therefore r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|} \quad (7.8)$$

## 7.1.2 Multiple classes

Consider the extension of linear discriminants to  $K > 2$  classes, building a  $K$ -class discriminants by combining a number of two-class discriminant functions on which leads to some serious difficulties (Duda and Hart, 1973).

Consider the use of  $K - 1$  classifier each of which solves a two-class problem of separating points into particular class  $\mathcal{C}_k$  from points not in that class. This is known as a **one-versus-the-rest** classifier. An alternative is to introduce  $K(K - 1)/2$  binary discriminant functions, one for every possible pair classes. This is known as a **one-versus-one** classifier. Both of these two approaches run into the problem of ambiguous regions. We can avoid these difficulties by considering a single  $K$ -class discriminant comprising  $k$  linear functions of the form

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (7.9)$$

and then assigning a point  $\mathbf{x}$  to class  $\mathcal{C}_k$  if  $y_k(\mathbf{x}) > y_j(\mathbf{x})$  for all  $j \neq k$ . The decision boundary between class  $\mathcal{C}_k$  and  $\mathcal{C}_j$  is therefore given by  $y_k(\mathbf{x}) = y_j(\mathbf{x})$  and hence corresponds to a  $(D - 1)$ -dimensional hyperplane defined by

$$(\mathbf{w}_k - \mathbf{w}_j)^T + (w_{k0} - w_{j0}) = 0 \quad (7.10)$$

The decision regions of such a discriminant are always singly connected and convex.

## 7.1.3 Least squares for classification

Even as a discriminant function (where we use it to make decisions directly and dispense with any probabilistic interpretation) it suffers from some severe problems. Least-squares solutions lack robustness to outliers, and this applies equally to the classification application. Additional points (outliers) produce a significant change in the location of the decision boundary. The sum-of-squares error function penalizes predictions that are 'too correct' in that they lie a long way on the correct side of the decision boundary. More than lack of robustness, least-squares solutions may give poor results on classification problems.

The failure of least squares lies in the fact that it corresponds to maximum likelihood under the assumption of a Gaussian conditional distribution, whereas binary target vectors clearly have a distribution that is far from Gaussian (a Bernoulli).

## 7.1.4 Fisher's linear discriminant

### 7.1.4.1 representation

One way to view a linear classification model is in terms of dimensionality reduction. Consider first the case of two classes, and suppose we take the  $D$ -dimensional input vector  $\mathbf{x}$  and project it down to one dimension using

$$y = \mathbf{w}^T \mathbf{x}. \quad (7.11)$$

If we place a threshold on  $y$  and classify  $y \geq -w_0$  as class  $\mathcal{C}_\infty$ , and otherwise class  $\mathcal{C}_\epsilon$ , then we obtain our standard linear classifier discussed in the previous section.

#### 7.1.4.2 evaluation

By adjusting the components of the weight vector  $\mathbf{w}$ , we can select a projection that maximizes the class separation. To begin with, consider a two-class problem in which there are  $N_1$  points of class  $\mathcal{C}_\infty$  and  $N_2$  points of class  $\mathcal{C}_\epsilon$ , so that the mean vectors of the two classes are given by

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_\infty} \mathbf{x}_n, \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_\epsilon} \mathbf{x}_n \quad (7.12)$$

The simplest measure of the separation of the classes, when projected onto  $\mathbf{w}$ , is the separation of the projected class means. This suggests that we might choose  $\mathbf{w}$  so as to maximize

$$m_2 - m_1 = \mathbf{x}^T (\mathbf{m}_2 - \mathbf{m}_1) \quad (7.13)$$

where

$$m_k = \mathbf{w}^T \mathbf{m}_k \quad (7.14)$$

is the mean of the projected data from class  $\mathcal{C}_\parallel$ . However this expression can be made arbitrarily large simply by increasing the magnitude of  $\mathbf{w}$ . To solve this, we could constrain  $\mathbf{w}$  to have unit length, so that  $\sum_i w_i^2 = 1$ . Using a Lagrange multiplier to perform the constrained maximization, we then find that  $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$ . There is still a problem with this approach that the projected data have considerable overlap for strongly nondiagonal covariances of the class distributions.

#### 7.1.4.3 optimization

The idea proposed by Fisher is to **maximize a function that will give a large separation between the projected class means while also giving a small variance within each class, thereby minimizing the class overlap**. The projection formula 7.11 transforms the set of labelled data points in  $\mathbf{x}$  into a labelled set in the one-dimensional space  $y$ . The within-class variance of the transformed data from class  $\mathcal{C}_\parallel$  is therefore given by

$$s_k^2 = \sum_{n \in \mathcal{C}_\parallel} (y_n - m_k)^2 \quad (7.15)$$

where

$$y_n = \mathbf{w}^T \mathbf{x}_n \quad (7.16)$$

We can define the total within-class variance for the whole data set to be simply  $s_1^2 + s_2^2$ . The Fisher criterion is defined to be the ratio of the between-class variance to the within-class variance and is given by

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \quad (7.17)$$

We can make the dependence on  $\mathbf{w}$  explicit by rewrite the Fisher criterion in the form

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad (7.18)$$

where  $\mathbf{S}_B$  is the *between-class* covariance matrix, given by

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \quad (7.19)$$

so

$$\mathbf{w}^T \mathbf{S}_B \mathbf{w} = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w} = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1) \quad (7.20)$$

and  $\mathbf{S}_W$  is the total *within-class* covariance matrix, given by

$$\mathbf{S}_W = \sum_C \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T \quad (7.21)$$

so

$$\mathbf{w}^T \mathbf{S}_W \mathbf{w} = \sum_{C_k} \sum_{n \in \mathcal{C}_k} \mathbf{w}^T (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T = \sum_{C_i} \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2 \quad (7.22)$$

Differentiating 7.18 with respect to  $\mathbf{w}$ , we find that  $J(\mathbf{w})$  is maximized when

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w} \quad (7.23)$$

Rewriting this as

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w} \quad (7.24)$$

where

$$\lambda = \frac{(\mathbf{w}^T \mathbf{S}_B \mathbf{w})}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})} \quad (7.25)$$

which is called a **generalized eigenvalue** problem solution.

*Proof.*

$$\therefore \begin{cases} \frac{\partial}{\partial x} \frac{f(x)}{g(x)} = \frac{f'g - fg'}{g^2}, \text{ where } f' = \frac{\partial f(x)}{\partial x} \text{ and } g' = \frac{\partial g(x)}{\partial x} \\ \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \end{cases} \quad (7.26)$$

$$\therefore \quad (7.27)$$

$$\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}) = \frac{2\mathbf{S}_B \mathbf{w} (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) - (\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w}}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} \quad (7.28)$$

Setting the derivative to zero, we can get the result.

In particular, since

$$\lambda \mathbf{w} = \mathbf{S}_W (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1) \quad (7.29)$$

Multiply both sides of 7.24 by  $\mathbf{S}_W^{-1}$  we then obtain

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1) \quad (7.30)$$

Note that if the within-class covariance is isotropic, so that  $\mathbf{S}_W$  is proportional to the unit matrix ( $\mathbf{S}_W \propto \mathbf{I}$ ), we find that  $\mathbf{w}$  is proportional to the difference vector of the class means.

The result 7.30 is known as *Fisher's linear discriminant*, although strictly it is not a discriminant but rather a specific choice of direction for projection of the data down to one dimension. However, the projected data can subsequently be used to construct a discriminant, by choosing a threshold  $y_0$  so that we classify a new point as belonging to  $\mathcal{C}_1$  if  $y(\mathbf{x}) \geq y_0$  and can classify it as belonging to  $\mathcal{C}_2$  otherwise.

## 7.1.5 Extensions to higher dimensions and multiple classes

We can extend the above idea to multiple classes, and to higher dimensional subspaces, by finding a projection *matrix*  $\mathbf{W}$  which maps from  $D$  to  $L$  so as to maximize

$$J(\mathbf{W}) = \frac{|\mathbf{W} \Sigma_B \mathbf{W}^T|}{|\mathbf{W} \Sigma_W \mathbf{W}^T|} \quad (7.31)$$

where

$$\Sigma_B \triangleq \sum_c \frac{N_c}{N} (\mathbf{m}_c - \mathbf{m})(\mathbf{m}_c - \mathbf{m})^T \quad (7.32)$$

$$\Sigma_W \triangleq \sum_c \frac{N_c}{N} \Sigma_c \quad (7.33)$$

$$\Sigma_c \triangleq \sum_{i:y_i=c} (\mathbf{x}_i - \mathbf{m}_c)(\mathbf{x}_i - \mathbf{m}_c)^T \quad (7.34)$$

The solution can be shown to be

$$\mathbf{W} = \Sigma_W^{-\frac{1}{2}} \mathbf{U} \quad (7.35)$$

where  $\mathbf{U}$  are the  $L$  leading eigenvectors of  $\Sigma_W^{-\frac{1}{2}} \Sigma_B \Sigma_W^{-\frac{1}{2}}$ , assuming  $\Sigma_W$  is non-singular. (If it is singular, we can first perform PCA on all the data).

## 7.1.6 Probabilistic interpretation of FLDA \*

### 7.1.7 Relation to least squares

### 7.1.8 Fisher's discriminant for multiple classes

### 7.1.9 The perceptron algorithm

Another example of a linear discriminant model is the perceptron of Rosenblatt(1962), which corresponds to a two-class transformation to give a feature vector  $\phi(\mathbf{x})$ , and this is then used to construct a generalized linear model of the form

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x})) \quad (7.36)$$

where the nonlinear activation function  $f(\cdot)$  is given by a step function of the form

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a \leq 0 \end{cases} \quad (7.37)$$

#### 7.1.9.1 error function

A natural choice of error function would be the total number of misclassified patterns, which does not lead to a simple learning algorithm because the error is a piecewise constant function of  $\mathbf{w}$ , with discontinuities wherever a change in  $\mathbf{w}$  causes the decision boundary to move across one of the data points. Methods based on gradient can't be applied. An alternative error function is known as the **perceptron criterion**, given by

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi_n(\mathbf{x}) t_n \quad (7.38)$$

where we use  $t \in \{-1, +1\}$  coding scheme,  $\mathcal{M}$  denotes the set of all misclassified patterns. The perceptron criterion associates zero error with any pattern that is correctly classified, whereas for a misclassified pattern  $\mathbf{x}$  it tries to minimize the quantity  $-\mathbf{w}^T \phi_n(\mathbf{x}) t_n$ . The perceptron learning rule is not guaranteed to reduce the total error function at each stage.

However, the **perceptron convergence theorem** states that if there exists an exact solution (linearly separable), then the perceptron learning algorithm is guaranteed to find an exact solution in a finite number of steps.

Aside from difficulties with the learning algorithm, the perceptron does not provide probabilistic outputs, nor does it generalize readily to  $K > 2$  classes. The most important limitation, however, arises from the fact that it is based on linear combinations of fixed basis functions.

### 7.1.9.2 optimization

We now apply the stochastic gradient descent algorithm to this error function. The change in the weight vector  $\mathbf{w}$  is then given by

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{\tau} + \eta \phi_n t_n \quad (7.39)$$

where  $\eta$  is the learning rate parameter and  $\tau$  is an integer that indexes the steps of the algorithm. Because the perceptron function  $y(\mathbf{x}, \mathbf{w})$  is unchanged if we multiply  $\mathbf{w}$  by a constant, so we can set the learning rate  $\eta$  equal to 1 without loss of generality.

## 7.2 Probabilistic Generative Models

We turn next to a probabilistic view of classification and show how models with linear decision boundaries arise from simple assumptions about the distribution of the data, adopting a generative approach in which we model the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$ , as well as the class priors  $p(\mathcal{C}_k)$ , and then use these to compute posterior probabilities  $p(\mathcal{C}_k|\mathbf{x})$  through Bayes' theorem.

### 7.2.1 Two classes

Consider first of all the case two classes. The posterior probability for class  $\mathcal{C}_1$  can be written as

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \quad (7.40)$$

$$= \frac{1}{1 + \exp(-a)} = \sigma(a) \quad (7.41)$$

where we have defined

$$a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \quad (7.42)$$

and  $\sigma(a)$  is the **logistic sigmoid** function defined by

$$\sigma(a) = \frac{1}{1 + \exp(-a)} \quad (7.43)$$

The term 'sigmoid' means S-shaped, sometimes called a 'squashing function' because it maps the whole real axis into a finite interval. It satisfies symmetry property

$$\sigma(-a) = 1 - \sigma(a) \quad (7.44)$$

The inverse of the logistic sigmoid is given by

$$a = \ln \left( \frac{\sigma}{1 - \sigma} \right) \quad (7.45)$$

and is known as the **logit** function. It represents the log of the ratio of probabilities  $\ln[p(\mathcal{C}_1|\mathbf{x})/p(\mathcal{C}_2|\mathbf{x})]$  for the two classes, also known as the **log odds**.

## 7.2.2 Multiple classes

For the case of  $K > 2$  classes, we have posterior

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} \quad (7.46)$$

$$= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad (7.47)$$

which is known as the **normalized exponential** and can be regarded as a multiclass generalization of the logistic sigmoid. Here the quantities  $a_k$  are defined by

$$a_k = \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k) \quad (7.48)$$

The normalized function is also known as the **softmax function**, as it represents a smoothed version of the 'max' function because, if  $a_k \gg a_j$  for all  $j \neq k$ , then  $p(\mathcal{C}_k|\mathbf{x}) \simeq 1$  and  $p(\mathcal{C}_j|\mathbf{x}) \simeq 0$ .

The following investigate the consequences of choosing specific forms for the class-conditional densities, continuous and discrete inputs.

## 7.2.3 Continuous inputs

Assume that the class-conditional densities are **Gaussian** sharing the **same covariance matrix** and then explore the resulting form for the posterior probabilities. Thus the density (pdf) for class  $\mathcal{C}_k$  is given by

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma^{-1}(\mathbf{x} - \mu_k)\right\} \quad (7.49)$$

In the case of two classes, we have

$$a(\mathbf{x}) = \log \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \quad (7.50)$$

$$= -\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1) + \frac{1}{2}(\mathbf{x} - \mu_2)^T \Sigma^{-1}(\mathbf{x} - \mu_2) + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)} \quad (7.51)$$

$$= (\mu_1 - \mu_2)^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \log \frac{p(\mathcal{C}_1)}{\mathcal{C}_2} \quad (7.52)$$

$$= \mathbf{w}^T \mathbf{x} + w_0 \quad (7.53)$$

where

$$\mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2) \quad (7.54)$$

$$w_0 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \log \frac{p(\mathcal{C}_1)}{\mathcal{C}_2} \quad (7.55)$$

We see that the quadratic terms in  $\mathbf{x}$  from the exponents of the Gaussian densities have cancelled (due to the common covariance matrices) leading to a linear function of  $\mathbf{x}$  in the argument of logistic sigmoid. The prior  $p(\mathcal{C}_k)$  enter only through the bias parameter  $w_0$  so that it have effect of the parallel contours of constant posterior probability.

For the general case of  $K$  cases we have

$$a_k(\mathbf{x}) = \log p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k) \quad (7.56)$$

$$= \mu_k^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \mathbf{x}^T \Sigma^{-1} \mathbf{x} + \log p(\mathcal{C}_k) + const \quad (7.57)$$

Notice that  $\mathbf{x}\Sigma^{-1}\mathbf{x}$  is cancelled in the exponent of the softmax function. Therefore we can write

$$a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (7.58)$$

where

$$\mathbf{w}_k = \Sigma^{-T} \boldsymbol{\mu}_k \quad (7.59)$$

$$w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(\mathcal{C}_k) \quad (7.60)$$

We see that the  $a_k(\mathbf{x})$  are again linear functions of  $\mathbf{x}$  due to the shared covariances, otherwise we obtain quadratic functions of  $\mathbf{x}$ , giving rise to a **quadratic discriminant**.

## 7.2.4 Maximum likelihood solution

Having specified a parametric functional form for the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$ , we can then determine the values of the parameters, together with the prior class probabilities  $p(\mathcal{C}_k)$ , using **maximum likelihood**.

Consider first the case of two classes, each having a Gaussian class-conditional density with a shared covariance matrix, and suppose we have a data set  $\{\mathbf{x}_n, t_n\}$  where  $n = 1, \dots, N$ . Here  $t_n = 1$  denotes class  $\mathcal{C}_1$  and 0 denotes  $\mathcal{C}_2$ . Denote the prior class probability  $p(\mathcal{C}_1) = \pi$  so that  $p(\mathcal{C}_2) = 1 - \pi$ . For a data point  $\mathbf{x}_n$  from class  $\mathcal{C}_1$ , we have  $t_n = 1$  and hence

$$p(\mathbf{x}_n, \mathcal{C}_1) = p(\mathcal{C}_1)p(\mathbf{x}|\mathcal{C}_1) = \pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \quad (7.61)$$

For class  $\mathcal{C}_2$ , and hence

$$p(\mathbf{x}_n, \mathcal{C}_2) = p(\mathcal{C}_2)p(\mathbf{x}|\mathcal{C}_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma) \quad (7.62)$$

Thus the likelihood is given by

$$p(\mathbf{t} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma)]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma)]^{1-t_n} \quad (7.63)$$

where  $\mathbf{t} = (t_1, \dots, t_N)^T$ . Maximize the log likelihood

$$\log p(\mathbf{t} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2) = \sum_{n=1}^N \{t_n \log \pi + t_n \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) + (1 - t_n) \log(1 - \pi) + (1 - t_n) \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma)\} \quad (7.64)$$

The terms that depend on  $\pi$  are

$$\sum_{n=1}^N \{t_n \log \pi + (1 - t_n) \ln(1 - \pi)\} \quad (7.65)$$

Setting the derivative with respect to  $\pi$  equal to zero and rearranging, we obtain

$$\sum_{n=1}^N \left\{ t_n \frac{1}{\pi} + (1 - t_n) \frac{1}{1 - \pi} \right\} = 0 \quad (7.66)$$

$$\Rightarrow \pi \sum_{n=1}^N (t_n - 1) - \pi \sum_{n=1}^N t_n + \sum_{n=1}^N t_n = 0 \quad (7.67)$$

$$\Rightarrow \pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2} \quad (7.68)$$

where  $N_k$  denotes the total number of data points in class  $\mathcal{C}_k$ . Thus the maximum likelihood estimate for  $\pi$  is simply the fraction of points in class as expected, which can be easily generalized to the multiclass case where ... .

Now consider the maximization with respect to  $\boldsymbol{\mu}_1$ . The log likelihood function those terms that depend on  $\boldsymbol{\mu}_1$



$$\sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) = -\frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) + \text{const} \quad (7.69)$$

Setting the derivative with respect to  $\boldsymbol{\mu}_1$  to zero and arranging ,we obtain

$$\boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n \quad (7.70)$$

which is simply the mean of all the input vectors  $\mathbf{x}_n$  assigned to corresponding class. Similarly,

$$\boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) \mathbf{x}_n \quad (7.71)$$

Finally,consider the maximum likelihood solution for the shared covariance matrix  $\boldsymbol{\Sigma}$ .We have

$$-\frac{1}{2} \sum_{n=1}^N t_n \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \quad (7.72)$$

$$-\frac{1}{2} \sum_{n=1}^N (1-t_n) \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (1-t_n) (\mathbf{x}_n - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \quad (7.73)$$

$$= -\frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{N}{2} \text{Tr}\{\boldsymbol{\Sigma}^{-1} \mathbf{S}\} \quad (7.74)$$

where we have defined

$$\mathbf{S} = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2 \quad (7.75)$$

$$\mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T \quad (7.76)$$

$$\mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T \quad (7.77)$$

Using the standard result for the maximum likelihood solution for a Gaussian distribution,we see that

$$\boldsymbol{\Sigma} = \mathbf{S} \quad (7.78)$$

which represents a weighted average of the covariance matrices associated with each of the two classes separately.<sup>20</sup>

This result is easily extended to the  $K$  class problem.Note that the approach of fitting Gaussian distributions to the classes is not robust to outliers,because the maximum likelihood estimation of a Gaussian is not robust.

## 7.2.5 Discrete features

Assume binary feature value  $x_i \in \{0,1\}$  is Bernoulli distributed,thus we have class-conditional distributions of the form

$$p(\mathbf{x} | \mathcal{C}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i} \quad (7.79)$$

which contains  $D$  independent parameters for each class,giving

$$a_k(\mathbf{x}) = \sum_{i=1}^D \{x_i \ln \mu_{ki} + (1 - x_i) \ln (1 - \mu_{ki})\} + \ln p(\mathcal{C}_k) \quad (7.80)$$

which again are linear functions of the input values  $x_i$ . Analogous results are obtained for discrete variables each of which can take  $M > 2$  states.

## 7.2.6 Exponential family

As we have seen so far, for both Gaussian distributed and discrete inputs, the posterior class probabilities are given by generalized linear models with logistic sigmoid ( $K = 2$  classes) or softmax ( $K \gg 2$  classes) activation functions. These are particular cases of a more general result obtained by assuming that the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$  are members of the exponential family of distributions. Using the form 2.6 for members of the exponential family, we see that the distribution of  $\mathbf{x}$  can be written in the form

$$p(\mathbf{x}|\boldsymbol{\lambda}_k) = h(\mathbf{x})g(\boldsymbol{\lambda}_k) \exp\{\boldsymbol{\lambda}_k^T \mathbf{u}(\mathbf{x})\} \quad (7.81)$$

We restrict attention to the subclass for which  $\mathbf{u}(\mathbf{x}) = \mathbf{x}$ . Making use of 2.62 to introduce a scaling parameter  $s$ , then

$$p(\mathbf{x}|\boldsymbol{\lambda}_k, s) = \frac{1}{s} h(\mathbf{x}/s) g(\boldsymbol{\lambda}_k) \exp\left\{\frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{u}(\mathbf{x})\right\} \quad (7.82)$$

Note that each class has its own parameter vector  $\boldsymbol{\lambda}_k$  but share the same scale parameter  $s$ .

For the two-class problem, the posterior is given by a logistic sigmoid acting on a linear function

$$a(\mathbf{x}) = (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_1) - \ln g(\boldsymbol{\lambda}_2) + \ln p(\mathcal{C}_1) - \ln p(\mathcal{C}_2) \quad (7.83)$$

Similarly, for the  $K$ -class problem,

$$a_k(\mathbf{x}) = \boldsymbol{\lambda}_k^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_k) + \ln p(\mathcal{C}_k) \quad (7.84)$$

and so again is a linear function of  $\mathbf{x}$ .

## 7.3 Probabilistic Discriminative Models

We have seen that the posterior probability for the two-class classification and multiclass case can be written as a logistic and softmax function respectively for a wide choice of class-conditional distributions  $p(\mathbf{x}|\mathcal{C}_k)$ . For a specific class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$ , we use maximum likelihood to determine the parameters of the densities as well as the class priors  $p(\mathcal{C}_k)$  and use Bayes' theorem to find the posterior class probabilities. However, an alternative approach is to use the functional form of the generalized linear model explicitly and to determine its parameters directly by using maximum likelihood, such as **iterative reweighted least squares** or IRLS.

The indirect approach to finding the parameters of a generalized linear model by fitting class-conditional densities and class priors separately and then applying Bayes' theorem, represents an example of **generative modelling** because we could take such a model and generate synthetic data by drawing values of  $\mathbf{x}$  from the marginal distribution  $p(\mathbf{x})$ . In the direct approach, we maximize a likelihood function defined through the conditional distribution  $p(\mathcal{C}_k|\mathbf{x})$ , which represents a form of **discriminative** training.

### 7.3.1 Fixed basis functions

So far, we have considered classification models that work directly with the original input vector  $\mathbf{x}$ . However, all of the algorithms can equally be applicable if we first make a fixed nonlinear transformation of the inputs using a vector of basis functions  $\phi(\mathbf{x})$ . The resulting decision boundaries will be linear in the feature space  $\phi$ , corresponding to nonlinear decision boundaries in the original  $\mathbf{x}$  space. One of the basis functions is typically set to a constant, say  $\phi_0(\mathbf{x}) = 1$ , so that the corresponding parameter  $w_0$  plays the role of a bias.

Note that nonlinear transformations cannot remove class overlap between the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$ .

## 7.3.2 Logistic regression

### 7.3.2.1 representation

We begin with two-class classification. The posterior probability of class  $C_1$  can be written as a logistic sigmoid acting on a linear function of the feature vector  $\phi$  so that

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad (7.85)$$

which has a linear dependence on the number of parameters. Logistic regression model can also be written as

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Ber}(y|\text{sigm}(\mathbf{w}^T \phi)) \quad (7.86)$$

where  $\mathbf{w}$  and  $\phi$  are extended vectors, i.e.,  $\mathbf{w} = (b, w_1, w_2, \dots, w_D)$ ,  $\phi = (1, \phi_1, \phi_2, \dots, \phi_D)$ .

Note that the derivative of the logistic sigmoid function is

$$\frac{d\sigma}{da} = \sigma(1 - \sigma) \quad (7.87)$$

### 7.3.2.2 evaluation

For a data set  $\{\phi_n, t_n\}$ , where  $t_n \in \{0, 1\}$  and  $\phi_n = \phi(\mathbf{x}_n)$ , with  $n = 1, \dots, N$ , the likelihood function can be written

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad (7.88)$$

where  $\mathbf{t} = (t_1, \dots, t_N)^T$  and  $y_n = p(C_1|\phi_n)$ . The negative logarithm of the likelihood, which gives the **cross-entropy** error function is in the form

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (7.89)$$

Taking the gradient of the error function with respect to  $\mathbf{w}$ , we obtain

$$\nabla E(\mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \frac{y_n(1 - y_n)}{y_n} - (1 - t_n) \frac{y_n(1 - y_n)}{1 - y_n} \right\} \phi_n \quad (7.90)$$

$$= -\sum_{n=1}^N \{t_n(1 - y_n) - (1 - t_n)y_n\} \phi_n \quad (7.91)$$

$$= \sum_{n=1}^N (y_n - t_n) \phi_n \quad (7.92)$$

We see that the factor involving the derivative of the logistic sigmoid has cancelled. In particular, the contribution to the gradient from data point  $n$  is given by the 'error'  $y_n - t_n$  between the target value and the prediction of the model, times the basis function vector  $\phi_n$ , taking the same form as the gradient of the sum-of-squares error function for the linear regression model.

Note that maximum likelihood can exhibit severe **over-fitting** for data sets that are linearly separable. The singularity can be avoided by inclusion of a prior and finding a **MAP(maximum a posterior)** solution for  $\mathbf{w}$ , or equivalently by adding a **regularization** term to the error function.

### 7.3.2.3 optimization

See the following contents.

### 7.3.3 Iterative reweighted least squares

#### 7.3.3.1 Newton-Raphson update

For logistic regression, there is no longer a closed-form solution, due to the nonlinearity of the logistic sigmoid function. The error function is concave, and can be minimized by an efficient iterative technique based on the **Newton-Raphson** iterative optimization scheme, which uses a local quadratic approximation to the log likelihood function.

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} - \mathbf{H}^{-1} \nabla E(\mathbf{w}) \quad (7.93)$$

where  $\mathbf{H}$  is the Hessian matrix whose elements comprise the second derivatives of  $E(\mathbf{w})$  with respect to the components of  $\mathbf{w}$ .

Now apply the Newton-Raphson method to the linear regression model with sum-of-squares error function. The gradient and Hessian of this error function are given by

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \phi_n - t_n) \phi_n = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t} \quad (7.94)$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \phi_n \phi_n^T = \Phi^T \Phi \quad (7.95)$$

where  $\Phi$  is the  $N \times N$  design matrix, whose  $n^{th}$  row is given by  $\phi_n^T$ . The Newton-Raphson update then takes the form

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} - (\Phi^T \Phi)^{-1} \{ \Phi^T \Phi \mathbf{w}^{(old)} - \Phi^T \mathbf{t} \} \quad (7.96)$$

$$= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (7.97)$$

which we recognize as the standard least-squares solution.

Now apply the Newton-Raphson update to the cross-entropy error function for logistic regression model. Gradient and Hessian of this error function are given by

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (\mathbf{y} - \mathbf{t}) \quad (7.98)$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi \quad (7.99)$$

where we have introduced the  $N \times N$  diagonal matrix  $\mathbf{R}$  with elements

$$R_{nn} = y_n (1 - y_n) \quad (7.100)$$

We see that the Hessian depends on  $\mathbf{w}$  and is positive definite.

The Newton-Raphson update formula for the logistic regression model then becomes

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} - (\Phi^T \mathbf{R} \Phi)^{-1} \quad (7.101)$$

$$= (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(old)} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \quad (7.102)$$

$$= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z} \quad (7.103)$$

where  $\mathbf{z}$  is an  $N$ -dimensional vector with elements

$$\mathbf{z} = \Phi \mathbf{w}^{(old)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}) \quad (7.104)$$

We see that the update formula takes the form of a set of normal equations for a weighted least-squares problem. Because the weighting matrix  $\mathbf{R}$  depends on the parameter vector  $\mathbf{w}$ , we must apply the normal equations iteratively, each time using the new weight vector  $\mathbf{w}$  to compute a revised weighting matrix  $\mathbf{R}$ . For this reason, the algorithm is known as **iterative reweighted least squares**, or IRLS. The elements of the diagonal weighting matrix  $\mathbf{R}$  can be

interpreted as variances because the mean and variance of  $t$  in the logistic regression model are given by

$$\mathbb{E}[t] = \sigma(\mathbf{x}) = y \quad (7.105)$$

$$\text{var}[t] = \mathbb{E}[t^2] - \mathbb{E}[t]^2 = \sigma(\mathbf{x}) - \sigma(\mathbf{x})^2 = y(1 - y) \quad (7.106)$$

where we have used the property  $t^2 = t$  for  $t \in \{0, 1\}$ .

### 7.3.3.2 Linearization

In fact, we can interpret IRLS as the solution to a linearized problem in the space of the variable  $a = \mathbf{w}^T \phi$ . Apply Taylor expansion and then we obtain

$$a_n(\mathbf{w}) \simeq a_n(\mathbf{w}^{(old)}) + \left. \frac{da_n}{dy_n} \right|_{\mathbf{w}^{(old)}} (t_n - y_n) \quad (7.107)$$

$$= \phi_n^T \mathbf{w}^{(old)} - \frac{(y_n - t_n)}{y_n(1 - y_n)} = z_n \quad (7.108)$$

## 7.3.4 Multiclass logistic regression

In our discussion of generative models for multiclass classification, the posterior probabilities are given by a softmax transformation of linear functions of the feature variables, so that

$$p(C_k | \phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad (7.109)$$

$$= \frac{1}{1 + \sum_{j \neq k} \exp(a_j - a_k)} \quad (7.110)$$

where the 'activations'  $a_k$  are given by

$$a_k = \mathbf{w}_k^T \phi \quad (7.111)$$

Here we consider the use of maximum likelihood to determine the parameters  $\{\mathbf{w}_k\}$  of this model directly. The derivatives of  $y_k$  with respect to all of the activations  $a_j$  are given by

$$\frac{\partial y_k}{\partial a_j} = \frac{\partial}{\partial a_j} \frac{1}{1 + \sum_{i \neq k} \exp(a_i - a_k)} \quad (7.112)$$

$$= - \frac{e^{a_i - a_k}}{(1 + \sum_{i \neq k} e^{a_i - a_k})^2} \quad (7.113)$$

$$= y_k(I_{kj} - y_j) \quad (7.114)$$

where  $I$  are the elements of the identity matrix.

Next we write down the likelihood function, using the 1-of- $K$  coding scheme.

$$p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k | \phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}} \quad (7.115)$$

where  $y_{nk} = y_k(\phi_n)$ , and  $\mathbf{T}$  is an  $N \times K$  matrix of target variables with elements  $t_{nk}$ . Taking the negative logarithm then gives

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \quad (7.116)$$

which is known as the *cross – entropy* error function for the multiclass classification problem.

Take the gradient of the error function, making use of the derivatives of the softmax function, we obtain

$$\nabla_{w_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n \quad (7.117)$$

To find a batch algorithm, we evaluate the Hessian matrix

$$\nabla_{w_k} \nabla_{w_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi_n \phi_n^T \quad (7.118)$$

The Hessian matrix is also positive definite and so the error function has a unique minimum.

### 7.3.5 Probit regression

A broad range of class-conditional distributions, described by the exponential family, the resulting posterior class probabilities are given by a logistic (or softmax) transformation acting on a linear function of the feature variables, but not all choices of class-conditional density give rise to such a simple form for the posterior probabilities.

The cumulative distribution function is given by

$$\Phi(a) = \int_{-\infty}^a \mathcal{N}(\theta|0, 1) d\theta \quad (7.119)$$

which is known as the **probit function**. It has a sigmoidal shape. A evaluation of a closely related function is

$$\text{erf}(a) = \frac{2}{\sqrt{\pi}} \int_0^a \exp(-\theta^2/2) d\theta \quad (7.120)$$

known as the **erf function** or **error function**. The generalized linear model based on a probit function is known as **probit regression**.

One issue that occurs in practical applications is that of **outliers**. Because probit activation function they decay like  $\exp(-x^2)$ , so probit model can be significantly more sensitive to outliers. However, the effect of mislabelling is easily incorporated into a probabilistic model by introducing a probability  $\varepsilon$  that the target value  $t$  has been flipped to the wrong value,

$$p(t|\mathbf{x}) = (1 - \varepsilon) \sigma(\mathbf{x}) + \varepsilon (1 - \sigma(\mathbf{x})) \quad (7.121)$$

$$= \varepsilon + (1 - 2\varepsilon) \sigma(\mathbf{x}) \quad (7.122)$$

where  $\sigma(\mathbf{x})$  is the activation function with input vector  $\mathbf{x}$ . Here  $\varepsilon$  may be set in advance, or it may be treated as a **hyperparameter** whose value is inferred from data.

### 7.3.6 Canonical link functions

We now show that there is a general result of assuming a conditional distribution for the target variable from the exponential family, along with a corresponding choice for the activation function known as the **canonical link function**.

Making use of the restricted form of exponential family distribution 7.82 for target variable  $t$

$$p(t|\eta, s) = \frac{1}{s} h\left(\frac{t}{s}\right) g(\eta) \exp\left\{\frac{\eta t}{s}\right\} \quad (7.123)$$

The conditional mean of  $t$ , which denoted by  $y$ , is given by

$$y \equiv \mathbb{E}[t|\eta] = -s \frac{d}{d\eta} \ln g(\eta) \quad (7.124)$$

Thus  $y$  and  $\eta$  must be related, and we denote this relation through  $\eta = \psi(y)$ .

We define a **generalized linear model** to be one for which  $y$  is a nonlinear function of a linear combination of the input(or feature) variables so that

$$y = f(\mathbf{w}^T \phi) \quad (7.125)$$

where  $f(\cdot)$  is known as the **activation function**, and  $f^{-1}(\cdot)$  is known as the **link function** in statistics.

The log likelihood function for this model, which, as a function of  $\eta$ , is given by

$$\ln p(\mathbf{t}|\eta, s) = \sum_{n=1}^N \ln p(t_n|\eta, s) = \sum_{n=1}^N \left\{ \ln g(\eta_n) + \frac{\eta_n t_n}{s} \right\} + \text{const} \quad (7.126)$$

where we assume that all observations share a common scale parameter (which corresponds to the noise variance for a Gaussian distribution for instance) and so  $s$  is independent of  $\eta$ . The derivative with respect to the model parameter  $\mathbf{w}$  is then given by

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\eta, s) = \sum_{n=1}^N \left\{ \frac{d}{d\eta_n} \ln g(\eta_n) + \frac{t_n}{s} \right\} \frac{d\eta_n}{dy_n} \frac{dy_n}{da_n} \nabla a_n \quad (7.127)$$

$$= \sum_{n=1}^N \frac{1}{s} \{t_n - y_n\} \psi'(y_n) f'(a_n) \phi_n \quad (7.128)$$

where  $a_n = \mathbf{w}^T \phi_n$ , and we have used  $y_n = f(a_n)$  together. There is a considerable simplification if we choose a particular form of the link function  $f^{-1}(y)$  given by

$$f^{-1}(y) = \psi(y) \quad (7.129)$$

In this case, the gradient of error function reduces to

$$\nabla \ln E(\mathbf{w}) = \frac{1}{s} \sum_{n=1}^N \{y_n - t_n\} \phi_n \quad (7.130)$$

For the Gaussian  $s = \beta^{-1}$ , whereas for logistic model  $s = 1$ .

### 7.3.7 The Laplace Approximation

### 7.3.8 Model comparison and BIC

## 7.4 Bayesian Logistic Regression

### 7.4.1 Laplace approximation

### 7.4.2 Predictive distribution





## Chapter 8

# Neural Networks

Applicability of models that comprised linear combinations of **fixed basis functions** is limited by the curse of dimensionality. In order to apply such models to large-scale problems, it is necessary to adapt the basis functions to the data.

Support vector machines and relevance vector machine **choose a subset from a fixed set of basis functions** and results in much sparser models.

An alternative approach is to fix the number of basis functions in advance but allow them to be **adaptive**, in other words to use parametric forms for the basis function

### 8.1 Feed-forward Network Functions

The linear models are based on linear combinations of fixed nonlinear basis functions  $\phi_j(\mathbf{x})$  and take the form

$$y(\mathbf{x}, \mathbf{w}) = f\left(\sum_{j=1}^M w_j \phi_j(\mathbf{x})\right) = f(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})) \quad (8.1)$$

where  $f(\cdot)$  is a nonlinear activation function in the case of classification and is the identity in the case of regression.

Neural network model can be described a series of functional transformations. First construct  $M$  linear combinations of the input variables  $x_1, \dots, x_D$  in the form

$$a_j = \sum_{i=1}^D w_{ji}^{(l)} + w_{j0}^{(l)} \quad (8.2)$$

where  $j = 1, \dots, M$ , and the superscript  $(l)$  indicates the  $l$ -th layer of the network. We refer the parameters  $w_{ji}^{(l)}$  as **weights** and  $w_{j0}^{(l)}$  as **biases**. The quantity  $a_j$  are known as **activations**. Each of them is then transformed using **activation function**  $h(\cdot)$  to give

$$z_j = h(a_j) \quad (8.3)$$

These quantities correspond to the outputs of the basis functions in linear model that, in the context of neural networks, are called **hidden units**. The nonlinear functions  $h(\cdot)$  are generally chosen to be sigmoidal or the ' $\tanh$ ' function.

The transformation of the following layer of the network, combine these values to give **output unit activations**

$$a_k = \sum_{j=1}^M w_{kj}^{(2)} z_j + w_{k0}^{(2)} \quad (8.4)$$

where  $k = 1, \dots, K$ , and  $K$  is the total number of outputs.

For regression,  $y_k = a_k$  and binary classification, we can use  $y_k = \sigma(a_k) = \frac{1}{1 + \exp(-a)}$ . Finally, for multiclass problems, a **softmax** activation function is used. **Unit activation function** is discussed later.

The process of evaluating the output of the network diagram is interpreted as a **forward propagation** of information through the network. Deterministic rather than stochastic variables of internal nodes do not represent probabilistic graphical models. Each (hidden or output) unit in such network computes a function given by

$$z_k = h\left(\sum_j w_{kj} z_j\right) \quad (8.5)$$

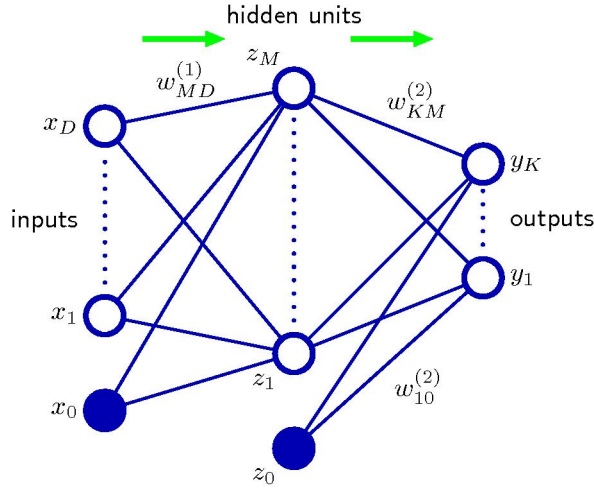


Fig. 8.1: Network diagram for two-layer neural network

where the sum runs over all units that send connections to unit  $k$  (including the bias parameter).

### 8.1.1 Weight-space symmetries

Consider a two-layer network as shown before. For  $M$  hidden units, there will be  $M$  'sign-flip' symmetries, i.e.,  $2^M$  equivalent weight vectors. Similarly, interchanging the values of all the weights (and the bias) can give rise to the same mapping function from inputs to outputs represented by the network. There will be  $M!$  equivalent weight vectors. The network has an overall weight-space symmetry factor of  $M!2^M$ .

## 8.2 Network Training

Given a training set comprising a set of input vectors  $\mathbf{x}_n$ , where  $n = 1, \dots, N$ , together with a corresponding set of target vectors  $\mathbf{t}_n$ , we minimize the error function. We can provide a probabilistic interpretation to the network outputs.

For regression, and for the moment we consider a single target variable  $t$  that can take any real value. We assume that  $t$  has a Gaussian distribution with an  $\mathbf{x}$ -dependent mean, which is given by the output of the neural network, so that

$$p(t|\mathbf{w}, \mathbf{w}) = \mathcal{N}(t|\mathbf{y}(\mathbf{x}, \mathbf{w}), \beta^{-1}) \quad (8.6)$$

where  $\beta$  is the precision (inverse variance) of the Gaussian noise. Maximize the likelihood function as in linear models. View the network as having an output activation function that is the identity, so that  $y_k = a_k$ . The corresponding sum-of-squares error function has the property

$$\frac{\partial E}{\partial a_k} = \frac{\partial E}{\partial y_k} = y_k - t_k \quad (8.7)$$

For binary classification in which  $t \in \{0, 1\}$ , consider a network having a single output whose activation function is a logistic sigmoid

$$y = \sigma(a) = \frac{1}{1 + \exp(-a)} \quad (8.8)$$

We can interpret  $y(\mathbf{x}, \mathbf{w})$  as the conditional probability  $p(C_1 | \mathbf{x})$ . the conditional distribution of targets given inputs is then a Bernoulli distribution of the form

$$p(t | \mathbf{x}, \mathbf{w}) = y(\mathbf{x}, \mathbf{w})^t (1 - y(\mathbf{x}, \mathbf{w}))^{1-t} \quad (8.9)$$

For  $K$  separate binary classifications, the conditional distribution of the targets is

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}) = \prod_{k=1}^K y_k(\mathbf{x}, \mathbf{w})^{t_k} (1 - y_k(\mathbf{x}, \mathbf{w}))^{1-t_k} \quad (8.10)$$

Taking the negative logarithm of the corresponding likelihood function, we get **cross-entropy** error function

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} + (1 - t_{nk}) \ln(1 - y_{nk}) \quad (8.11)$$

The derivative of the error function with respect to the **activation** for a particular output unit is

$$\frac{\partial E}{\partial a_k} = \frac{\partial - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} + (1 - t_{nk}) \ln(1 - y_{nk})}{\partial a_k} \quad (8.12)$$

$$= -t_k \frac{y_k(1 - y_k)}{y_k} + (1 - t_k) \frac{-y_k(1 - y_k)}{1 - y_k} \quad (8.13)$$

$$= y_k - t_k \quad (8.14)$$

Following the logistic regression, the **output unit activation** is given by the **softmax** function

$$y_k(\mathbf{x}, \mathbf{w}) = \frac{\exp(a_k(\mathbf{x}, \mathbf{w}))}{\sum_j \exp(a_j(\mathbf{x}, \mathbf{w}))} \quad (8.15)$$

### 8.2.1 Parameter optimization

Turn next to the task of finding a weight vector  $\mathbf{w}$  which minimizes the chosen function  $E(\mathbf{w})$ . Resort to **iterative**

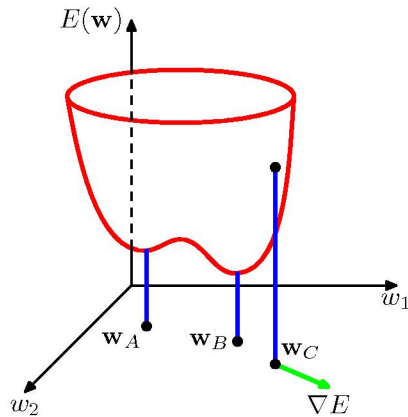


Fig. 8.2: geometrical view of the error function  $E(\mathbf{w})$  as a surface sitting over weight space. A is a local minimum and B is the global minimum. The local gradient of the error surface is given by the vector  $\nabla E$ .

numerical procedures. Most techniques involve choosing some initial value  $\mathbf{w}^{(0)}$  for weight vector and then moving through weight space in a succession of steps of the form

$$\mathbf{w}^{(\eta+1)} = \mathbf{w}^{(\eta)} + \Delta \mathbf{w}^{(\eta)} \quad (8.16)$$

Many algorithms make use of gradient information, evaluating  $\nabla E(\mathbf{w})$  at the new weight vector  $\mathbf{w}^{(\eta+1)}$ .

## 8.2.2 Local quadratic approximation

Consider the **Taylor expansion** of  $E(\mathbf{w})$  around some point in weight space.

## 8.2.3 Use of gradient information

It is possible to evaluate the gradient of an error function efficiently by means of the backpropagation procedure. Because each evaluation of  $\nabla E$  brings  $W$  items of information, we might hope to find the minimum of the function in  $O(W)$  gradient evaluations. Each such evaluation takes only  $O(W)$  steps and so the minimum can be now found in  $O(W^2)$  steps.

## 8.2.4 Gradient descent optimization

### Batch gradient descent

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w})^{(\tau)} \quad (8.17)$$

where the parameter  $\eta > 0$  is known as the **learning rate**.

More efficient methods than this poor algorithm are **conjugate gradients** and **quasi-Newton** methods, which are more robust and much faster. **On-line** version of gradient has also proved useful in practice. On-line gradient descent, known as **sequential gradient descent** or **stochastic gradient descent**, makes an update to the weight vector based on one data point at a time, so that

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)}) \quad (8.18)$$

## 8.3 Error Backpropagation

The goal in this section is to find an efficient technique for evaluating the gradient of an error function  $E(\mathbf{w})$  for a feed-forward neural network. This can be achieved using a local message passing scheme in which information is sent alternately forwards and backwards through the network, known as **error backpropagation** or **backprop**.

### 8.3.1 Evaluation of error-function derivatives

Many error functions of practical interest, for instance those defined by maximum likelihood for a set of i.i.d data, comprise a sum of terms, one for each data point, so that

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) \quad (8.19)$$

In a general feed-forward network, each unit computes a weighted sum of its inputs of the form

$$a_j = \sum_i w_{ji} z_i \quad (8.20)$$

where  $z_i$  is the activation of a unit that sends a connection to unit  $j$  and  $w_{ji}$  is the weight associated with that connection. The sum is transformed by an activation function  $h(\cdot)$  to give the activation  $z_j$  of unit  $j$  in the form

$$z_j = h(a_j) \quad (8.21)$$

Evaluate the derivate of  $E_n$  with respect to a weight  $w_{ji}$ , applying the **chain rule** for partial derivatives

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} \quad (8.22)$$

Now introduce a notation

$$\delta_j \equiv \frac{\partial E_n}{\partial a_j} \quad (8.23)$$

referred as **errors**. Using the sum form

$$\frac{\partial a_j}{\partial w_{ji}} = z_i \quad (8.24)$$

Substituting them back

$$\frac{\partial E_n}{\partial w_{ji}} = \delta_j z_i. \quad (8.25)$$

For the output units, we have

$$\delta_k = y_k - t_k \quad (8.26)$$

For hidden units, make use of the chain rule for partial derivatives

$$\delta_j \equiv \frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} \quad (8.27)$$

$$(8.28)$$

where the sum runs over all units  $k$  to which unit  $j$  sends connections.

$$\therefore \delta_k = \frac{\partial E_n}{\partial a_k} \quad (8.29)$$

$$\therefore \delta_j \equiv \frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} \quad (8.30)$$

$$= \sum_k \delta_k \frac{\partial a_k}{\partial a_j} \quad (8.31)$$

$$= \sum_k \delta_k h'(a_j) w_{kj} \quad (8.32)$$

$$= h'(a_j) \sum_k w_{kj} \delta_k \quad (8.33)$$

### 8.3.2 Efficiency of back propagation

### 8.3.3 The Jacobian matrix

Backpropagation can also be applied to the calculation of other derivatives. Consider the evaluation of the **Jacobian matrix**, whose elements are given by the derivatives of the network outputs with respects to the inputs

$$J_{ki} = \frac{\partial y_k}{\partial x_i} \quad (8.34)$$

## 8.4 The Hessian Matrix

## 8.5 Regularization in Neural Networks

## 8.6 Mixture Density Networks

## 8.7 Bayesian Neural Networks

## Chapter 9

# Sparse Kernel Machines

### 9.1 Introduction

One of significant limitations of kernel methods is that the kernel function  $\|(\mathbf{x}_n, \mathbf{x}_m)$  must be evaluated for all possible pairs  $\mathbf{x}_n$  and  $\mathbf{x}_m$  of training points, computationally infeasible. Kernel-based algorithms that have **sparse** solutions predict for new inputs depend only on the kernel function evaluated at a subset of the training data points, such as **support vector machine** (SVM).

### 9.2 Maximum Margin Classifiers

Two-class classification problem using linear models

$$y(x) = \mathbf{w}^T \phi(x) + b \quad (9.1)$$

where  $\phi(x)$  denotes a fixed feature-space transformation, and  $b$  is bias parameter.

For linear separable feature space, the parameters satisfies  $y(\mathbf{x}_n) > 0$  for points having  $t_n = +1$  and  $y(\mathbf{x}_n) < 0$  for  $t_n = -1$ , so that  $t_n y(\mathbf{x}_n) > 0$  for all points.

**Margin** is the smallest distance between the decision boundary and any of the samples. The maximum margin solution can be motivated using **computational learning theory**, also known as **statistical learning theory**.

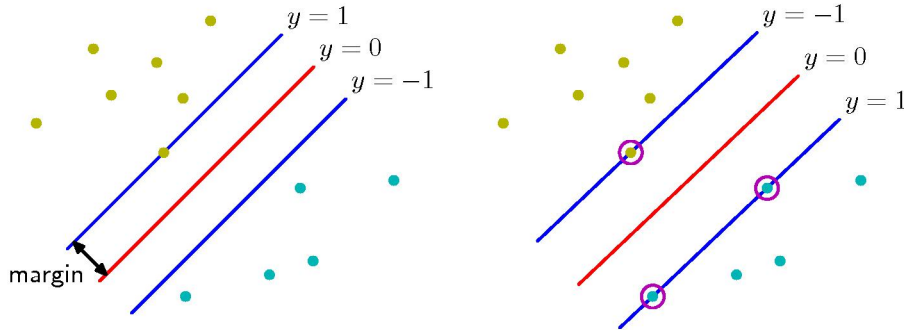


Fig. 9.1: margin

The **functional margin**  $t_n y(\mathbf{x}) > 0$  for data points correctly classified. The maximize it

$$\mathbf{w}, b = \arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\} \quad (9.2)$$

We can rescale parameters to set

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1 \quad (9.3)$$

for point closet to the surface. Then all data points satisfies the constraint

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, n = 1, \dots, N. \quad (9.4)$$

This is the **canonical representation of the decision hyperplane**. The optimization problem is equivalent to

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad (9.5)$$

subject to constraints 9.4, which is a **quadratic programming** problem.

Introducing Lagrange multipliers  $a_n \geq 0$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1\} \quad (9.6)$$

where  $\mathbf{a} = (a_1, \dots, a_N)^T$ . Setting the derivatives of  $L$  with respect to  $\mathbf{w}$  and  $b$  equal to zero, we obtain conditions

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \quad (9.7)$$

$$0 = \sum_{n=1}^N a_n t_n \quad (9.8)$$

Eliminating  $\mathbf{w}$  and  $b$  gives the **dual representation** of the maximum margin problem

$$\hat{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \quad (9.9)$$

with respect to  $\mathbf{a}$  subject to the constraints

$$a_n \geq 0, n = 1, \dots, N \quad (9.10)$$

$$\sum_{n=1}^N a_n t_n = 0. \quad (9.11)$$

The kernel function is defined by  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$ .

The solution to a quadratic programming problem in  $M$  variables has computational complexity of  $O(M^3)$ .

$y(\mathbf{x})$  can be expressed by

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b. \quad (9.12)$$

The *Karush – Kuhn – Tucker (KKT)* conditions require following the properties hold

$$a_n \geq 0 \quad (9.13)$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0 \quad (9.14)$$

$$a_n \{t_n y(\mathbf{x}_n) - 1 = 0\} \quad (9.15)$$

Data points for which  $a_n = 0$  will disappear and remaining ones are called **support vectors**. They lie on the maximum margin hyperplanes in feature space.

Having solved the quadratic programming problem, the threshold parameter  $b$

$$t_n \left( \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = 1 \quad (9.16)$$

where  $\mathcal{S}$  denotes the set of indices of the support vectors. Multiply through by  $t_n$ , making use of  $t_n^2 = 1$ , and then average these equations over all support vectors

$$b = \frac{1}{N_{\mathcal{S}}} \sum_{n \in \mathcal{S}} (t_n - \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m)) \quad (9.17)$$



where  $N_S$  is the total number of support vectors. Express the maximum-margin classifier in terms of the minimization of an error function with a quadratic regularizer

$$\sum_{n=1}^N E_{\infty}(y(\mathbf{x}_n)t_n - 1) + \lambda \|\mathbf{w}\|^2 \quad (9.18)$$

where  $E_{\infty}(z)$  is zero if  $z \geq 0$  and  $\infty$  otherwise to ensure the margin constraint.

### 9.2.1 Overlapping class distributions

In practice, the class-conditional distributions may overlap, in which case exact separation of the training data can lead to poor generalization. Introduce **slack variables**,  $\xi_n \geq 0$  where  $n = 1, \dots, N$ , one for each training data points. We allow points on the 'wrong side' but with a penalty that increases of the distance from the boundary.

$$\xi_n = |t_n - y(\mathbf{x}_n)| \quad (9.19)$$

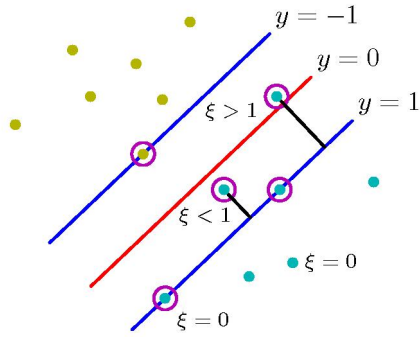


Fig. 9.2: slack variables  $\xi_n \geq 0$ . Data points with circles around them are support vectors

Then the classification constraint are replaced by

$$t_n y(\mathbf{x}_n) \geq 1 - \xi_n, n = 1, \dots, N \quad (9.20)$$

This is described as relaxing the hard margin constraint to give a **soft margin** and allows misclassification of training set data points.

Maximize the margin with softly penalized points on the wrong side of the margin boundary.

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2 \quad (9.21)$$

where the parameter  $C$  controls the trade-off between the slack variable penalty and the margin, minimizing errors and controlling model complexity. In the limit  $C \rightarrow \infty$ , the model gets more complex, less data points are misclassified.

Minimization with constraint

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n \quad (9.22)$$

where  $\{a_n \geq 0\}$  and  $\{\mu_n \geq 0\}$  are Lagrange multipliers. The KKT set of conditions are given by

$$a_n \geq 0 \quad (9.23)$$

$$t_n y(\mathbf{x}_n) - 1 + \xi_n \geq 0 \quad (9.24)$$

$$a_n(t_n y(\mathbf{x}_n) - 1 + \xi_n) = 0 \quad (9.25)$$

$$\mu_n \geq 0 \quad (9.26)$$

$$\xi_n \geq 0 \quad (9.27)$$

$$\mu_n \xi_n = 0 \quad (9.28)$$

where  $n = 1, \dots, N$ .

Optimize out  $\mathbf{w}, b$  and  $\{\xi_n\}$

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N a_n t_n \quad (9.29)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N a_n t_n = 0 \quad (9.30)$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n = C - \mu_n \quad (9.31)$$

Eliminated, the dual Lagrangian is in the form

$$\hat{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \|(x_n, x_m)\| \quad (9.32)$$

which is identical to the separable case, with different constraints:

$$0 \leq a_n \leq C \quad (9.33)$$

$$\sum_{n=1}^N a_n t_n = 0 \quad (9.34)$$

for  $n = 1, \dots, N$ , where the former are known as **box constraints**.

A subset of data points having  $a_n = 0$  do not contribute to the predictive model. Support vectors have  $a_n > 0$  and satisfy

$$t_n y(\mathbf{x}_n) = 1 - \xi_n \quad (9.35)$$

if  $a_n < C$ , then implies that  $\mu_n > 0$ , which requires  $\xi_n = 0$  and hence such points lie on the margin. Points with  $a_n = C$  can lie inside the margin and can either be correctly classified if  $\xi_n \leq 1$  or misclassified if  $\xi_n > 1$ .

To determine  $b$ , we note that support vectors for which  $0 \leq a_n \leq C$  have  $\xi_n = 0$  so that  $t_n y(\mathbf{x}_n) = 1$  and hence satisfy

$$t_n \left( \sum_{m \in \mathcal{S}} a_m t_m \|(x_n, x_m)\| + b \right) = 1 \quad (9.36)$$

A numerically stable solution is obtained by averaging:

$$b = \frac{1}{N_{\mathcal{M}}} \sum_{n \in \mathcal{M}} (t_n - \sum_{m \in \mathcal{S}} a_m t_m \|(x_n, x_m)\|) \quad (9.37)$$

where  $\mathcal{M}$  denotes the set of indices of data points having  $0 \leq a_n \leq C$ .

### 9.2.1.1 v-SVM

Maximize

$$\hat{L}(\mathbf{a}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \|(x_n, x_m) \quad (9.38)$$

subject to the constraints

$$0 \leq a_n \leq 1/N \quad (9.39)$$

$$\sum_{n=1}^N a_n t_n = 0 \quad (9.40)$$

$$\sum_{n=1}^N a_n \geq \nu \quad (9.41)$$

The parameter  $\nu$ , which replaces  $C$ , can be interpreted as both an upper bound on the fraction of **margin errors** and a lower bound on the fraction of support vectors.

### 9.2.1.2 Optimization

**chunking** Lagrangian is unchanged if we remove the rows and columns of the kernel matrix corresponding to Lagrange multipliers. Implemented using **protected conjugate gradients**.

**Decomposition methods** solves a series of smaller quadratic programming problems but are designed so that each of these is of a fixed size.

**sequential minimal optimization** or SMO Takes the concept of chunking to the extreme limit and considers just two Lagrange multipliers at a time. Heuristics are given for choosing the pair of Lagrange multipliers to be considered at each step.

Support vector machines don't manage to avoid the curse of dimensionality because there are constraints amongst the feature values that restrict the effective dimensionality of feature space.

### 9.2.1.3 Relation to Logistic Regression

The objective function can be written

$$\sum_{n=1}^N E_{SV}(y_n t_n) + \lambda \|\mathbf{w}\|^2 \quad (9.42)$$

where  $\lambda = (2C)^{-1}$ , and  $E_{SV}(\cdot)$  is the **hinge** error function defined by

$$E_{SV}(y_n t_n) = [1 - y_n t_n]_+ \quad (9.43)$$

where  $[\cdot]_+$  denotes the positive part.

For logistic regression, target variable  $t \in \{-1, 1\}$ , so

$$p(t|y) = \sigma(yt) \quad (9.44)$$

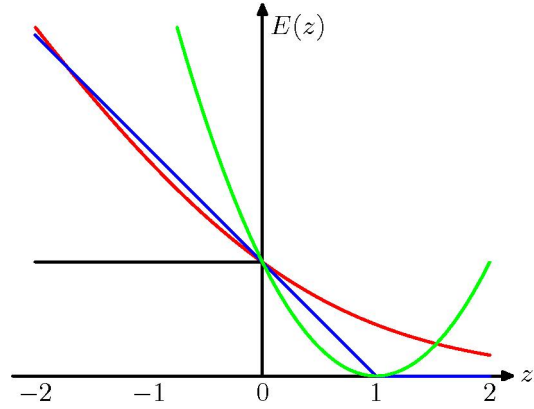
Construct an error function by taking the negative logarithm of likelihood function with a quadratic regularizer

$$\sum_{n=1}^N E_{LR}(y_n t_n) + \lambda \|\mathbf{w}\|^2 \quad (9.45)$$

where

$$E_{LR}(yt) = \ln(1 + \exp(-yt)) \quad (9.46)$$

Fig. 9.3: Plot of the hinge error function used in support vector machines, shown in blue, along with the error function for logistic regression, rescaled by a factor of  $1/\ln(2)$  so that it passes through the point  $(0, 1)$ , shown in red. Also shown are the misclassification error in black and the squared error in green.



#### 9.2.1.4 Multiclass SVMs

**one-versus-the-rest** approach:  $K$  separate SVMs for each class.

**one-versus-one**  $.K(K-1)/2$  different 2-class SVMs on possible pairs of classes, which can lead to ambiguities.

**single-class** support vector machines, which solve an unsupervised learning problem related to probability density estimation.

#### 9.2.1.5 SVMs for regression

In simple linear regression we minimize a regularized error function given by

$$\frac{1}{2} \sum_{n=1}^N \{y_n - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2. \quad (9.47)$$

To obtain **sparse solutions**, the quadratic error function is replaced by an  $\varepsilon$ -insensitive error function. For example

$$E_\varepsilon(y(\mathbf{x}) - t) = \begin{cases} 0, & \text{if } |y(\mathbf{x}) - t| < \varepsilon; \\ |y(\mathbf{x}) - t| - \varepsilon, & \text{otherwise} \end{cases} \quad (9.48)$$

We therefore minimize

$$C \sum_{n=1}^N E_\varepsilon(y(\mathbf{x}_n) - t_n) + \frac{1}{2} \|\mathbf{w}\|^2 \quad (9.49)$$

where  $y(\mathbf{x})$  is the prediction function. The (inverse) regularization parameter denoted  $C$ , appears in front of the error term.

Introducing two **slack variables** for each data point  $\mathbf{x}_n$ . The condition for target points to lie inside the  $\varepsilon$ -tube is that  $y_n - \varepsilon \leq t_n \leq y_n + \varepsilon$ , and for those outside:

$$\begin{cases} \xi_n & \geq 0 \\ \hat{\xi}_n & \geq 0, n = 1, \dots, N \end{cases} \quad (9.50)$$

$$y(\mathbf{x}_n) + \varepsilon < t_n \leq y(\mathbf{x}_n) + \varepsilon + \xi_n \quad (9.51)$$

$$y(\mathbf{x}_n) - \varepsilon > t_n \geq y(\mathbf{x}_n) - \varepsilon - \hat{\xi}_n \quad (9.52)$$

The error function for support vector regression can then be written as

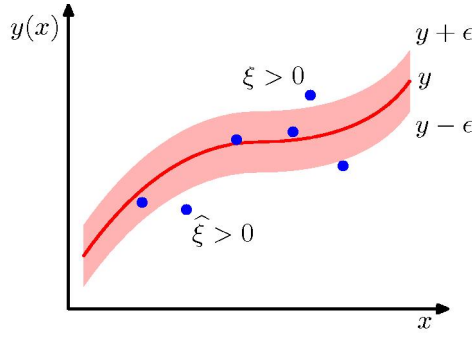


Fig. 9.4: Illustration of SVM regression, showing the regression curve together with the  $\epsilon$ -insensitive tube.

$$C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \| \mathbf{w} \|^2 \quad (9.53)$$

which must be minimized subject to the constraints. Introducing Lagrange multipliers

$$L = C \sum_{n=1}^N (\xi_n + \hat{\xi}_n) + \frac{1}{2} \| \mathbf{w} \|^2 - \sum_{n=1}^N (\mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n) \quad (9.54)$$

$$- \sum_{n=1}^N a_n (\epsilon + \xi_n + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (\epsilon + \hat{\xi}_n + y_n - t_n) \quad (9.55)$$

Substitute for  $y(x)$  using 9.1 and set the derivatives of the Lagrangian with respect to  $\mathbf{x}, b, \xi_n, \hat{\xi}_n$  to zero, giving

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N (a_n - \hat{a}_n) \phi(\mathbf{x}_n) \quad (9.56)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N (a_n - \hat{a}_n) = 0 \quad (9.57)$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n + \mu_n = 0 \quad (9.58)$$

$$\frac{\partial L}{\partial \hat{\xi}_n} = 0 \Rightarrow \hat{a}_n + \hat{\mu}_n = 0 \quad (9.59)$$

$$(9.60)$$

Using these to eliminate the corresponding variables, we see the dual problem of maximizing

$$\hat{L}(\mathbf{a}, \hat{\mathbf{a}}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) \|\mathbf{x}_n, \mathbf{x}_m\| \quad (9.61)$$

$$- \epsilon \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n \quad (9.62)$$

with respect to  $\{a_n\}, \{\hat{a}_n\}$ . We have the **box constraints**

$$0 \leq a_n \leq C \quad (9.63)$$

$$0 \leq \hat{a}_n \leq C \quad (9.64)$$

After substitution, the predictions can be made using

$$y(\mathbf{x}) = \sum_{n=1}^N (a_n - \hat{a}_n) \|\mathbf{x}, \mathbf{x}_n\| + b \quad (9.65)$$

The corresponding *karush – Kuhn – Tucker*(KKT) conditons for 9.54,which state that **at the solution the product of the dual variables and the constraints must vanish** are given by

$$a_n(\varepsilon + \xi_n + y_n - t_n) = 0 \quad (9.66)$$

$$\hat{a}_n(\varepsilon + \hat{\xi}_n - y_n + t_n) = 0 \quad (9.67)$$

$$(C - a_n)\xi_n = 0 \quad (9.68)$$

$$(C - \hat{a}_n)\hat{\xi}_n = 0 \quad (9.69)$$

The support vectors are those data points that contribute to predictions,in other words those for which either  $a_n \neq 0$  or  $\hat{a}_n \neq 0$ .These are points lying on the boundary of the  $\varepsilon$ -tube or outside the tube.

The parameter  $b$  satisfies

$$\varepsilon + y_n - t_n = 0 \quad (9.70)$$

for points  $0 < a_n < C$ .Solving for it

$$b = t_n - \varepsilon - \mathbf{w}^T \phi(\mathbf{x}_n) \quad (9.71)$$

$$= t_n - \varepsilon - \sum_{m=1}^N (a_m - \hat{a}_m) \|( \mathbf{x}_n, \mathbf{x}_m ) \quad (9.72)$$

An alternative formulation of SVM regression is  $\nu$  SVM.

$$\hat{L}(\mathbf{a}, \hat{\mathbf{a}}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) \|( \mathbf{x}_n, \mathbf{x}_m ) \quad (9.73)$$

$$- 0 \times \varepsilon \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n \quad (9.74)$$

subject to constraints

$$0 \leq a_n \leq C/N \quad (9.75)$$

$$0 \leq \hat{a}_n \leq C/N \quad (9.76)$$

$$\sum_{n=1}^N (a_n - \hat{a}_n) = 0 \quad (9.77)$$

$$\sum_{n=1}^N (a_n + \hat{a}_n) \leq \nu C \quad (9.78)$$

$$(9.79)$$

There are at most  $\nu N$  data points falling outside the insensitive tube,which at least  $\nu N$  data points are support vectors and so lie either on the tube or outside it.

### 9.3 Relevance Vector Machines

TODO

## Chapter 10

### Mixture models and EM

#### 10.1 Introduction

If we define a joint distribution over observed and latent variables, the corresponding distribution of the observed variables alone is marginalization. Mixture distributions can be interpreted in terms of discrete latent variables. As well as providing a framework for building more complex probability distributions, mixture models can also be used to cluster data.

#### 10.2 K-means Clustering

##### 10.2.1 representation

Suppose we have a data set  $\{x_1, \dots, x_N\}$  consisting of  $N$  observations of a random  $D$ -dimensional Euclidean variable  $x$ . Our goal is to partition the data set into some number  $K$  of clusters. First introduce a set of  $D$ -dimensional vectors  $\mu_k$ , where  $k = 1, \dots, K$ , in which  $\mu_k$  is a prototype associated with the  $k^{th}$  cluster, representing the centers of the clusters.

Our goal is then to find an assignment of data points to clusters, as well as a set of vectors  $\{\mu_k\}$ , such that the sum of squares of the distances of each data point to its closest vector  $\mu_k$  is a minimum.

A corresponding set of binary indicator variables  $r_{nk} \in \{0, 1\}$ , describing which of the  $K$  clusters the data point  $x_n$  is assigned to, so that if data point  $x_n$  is assigned to cluster  $k$  then  $r_{nk} = 1$ , and  $r_{nj} = 0$  for  $j \neq k$ . This is the 1-of- $K$  coding scheme.

##### 10.2.2 evaluation

Define the objective function, sometimes called a **distortion measure** as

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2 \quad (10.1)$$

##### 10.2.3 optimization

Our goal is to find values for the  $\{r_{nk}\}$  and the  $\{\mu_k\}$  so as to minimize  $J$ . We can do this through an iterative procedure in which each iteration involves two successive steps corresponding to successive optimizations with respect to the  $r_{nk}$  and the  $\mu_k$ .

Consider first the determination of the  $r_{nk}$ . We simply assign the  $n^{th}$  data point to the closest cluster centre.

$$r_{nk} = \begin{cases} 1, & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0, & \text{otherwise} \end{cases} \quad (10.2)$$

Now consider the optimization of the  $\mu_k$  with the  $r_{nk}$  held fixed. Setting the objective function  $J$ 's derivative with respect to  $\mu_k$  to zero giving:

$$\nabla_{\mu_k} J = 2 \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \mu_k) = 0 \quad (10.3)$$

$$\Rightarrow \mu_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}} \quad (10.4)$$

---

**Algorithm 1:** K-MEANS coordinate descent

---

**Input:** A set  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$

**Output:**  $\{r_{nk}\}$  and the  $\{\mu_k\}$  such that  $J$  is minimized.

1 initialize  $\mu_k$ ;

2 **repeat**

3     E(expectation).Minimize  $J$  with respect to the  $r_{nk}$ .  $r_{nk} = \begin{cases} 1, & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|^2 \\ 0, & \text{otherwise} \end{cases}$ , keeping  $\mu_k$  fixed

;

4     M(maximization).Minimize  $J$  with respect to the  $\mu_k$ :  $\mu_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$ , keeping  $r_{nk}$  fixed.;

5 **until** converged;

---

The objective may converge to a local rather than global minimum (MacQueen 1967).

Derive an on-line stochastic rather than the batch version of  $K$ -means, by applying Robbin-Monro procedure :

$$\mu_k^{new} = \mu_k^{old} + \eta_n (\mathbf{x}_n - \mu_k^{old}) \quad (10.5)$$

where  $\eta_n$  is the learning rate parameter, which is typically made to decrease monotonically as more data points are considered.

Generalize the  $K$ -means algorithm, introducing a general dissimilarity measure  $v(\mathbf{x}, \mathbf{x}')$

$$\hat{J} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} v(\mathbf{x}_n, \mu_k) \quad (10.6)$$

which gives the  **$K$ -medoids** algorithm. The computation cost of E step is  $O(KN)$ . The M step is potentially more complex. So it is common to restrict each cluster prototype to be equal to one of the data vectors assigned to that cluster, requiring  $O(N_k^2)$  evaluations of  $v(\cdot, \cdot)$ .

## 10.3 Mixtures of Gaussians

### 10.3.1 representation

The Gaussian mixture distribution can be written as a linear superposition of Gaussians in the form

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k) \quad (10.7)$$

Introduce a  $K$ -dimensional binary random variable  $\mathbf{z}$  having a 1-of- $K$  representation.

$$z_k \in \{0, 1\} \quad (10.8)$$

$$\sum_k z_k = 1 \quad (10.9)$$

The marginal distribution over  $\mathbf{z}$  is specified in terms of the mixing coefficients  $\pi_k$ , such that

$$p(z_k = 1) = \pi_k \quad (10.10)$$



where the parameters  $\{\pi_k\}$  must satisfy

$$0 \leq \pi_k \leq 1 \quad (10.11)$$

$$\sum_{k=1}^K \pi_k = 1 \quad (10.12)$$

Write this distribution in the form

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \quad (10.13)$$

The conditional distribution of  $\mathbf{x}$  given a particular value for  $\mathbf{z}$  is a Gaussian:

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (10.14)$$

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad (10.15)$$

The joint distribution is given by  $p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$ , and the marginal distribution of  $\mathbf{x}$  is obtained by summing the joint distribution over states of  $\mathbf{z}$ :

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (10.16)$$

which is a Gaussian mixture 10.7 .

The conditional probability of  $\mathbf{z}$  given  $\mathbf{x}$  can be evaluated by Bayes' theorem

$$\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \quad (10.17)$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (10.18)$$

We shall view  $\pi_k$  as the prior probability of  $z_k = 1$ , and the quantity  $\gamma(z_k)$  as the corresponding posterior.  $\gamma(z_k)$  can also be viewed as the **responsibility** that component  $k$  takes for 'explaining' the observation  $\mathbf{x}$ .

### 10.3.2 Maximum likelihood

Assuming that the data points are drawn independently from the same distribution (i.i.d), the log of the likelihood function is given by

$$\ln(p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})) = \sum_{n=1}^N \ln\left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \quad (10.19)$$

It is worth emphasizing that there are significant problems associated with the maximum likelihood framework applied to GMM, due to the presence of **singularities**. If one of the components of the mixture model 'collapse' onto a specific data point ( $\sigma_j \rightarrow 0$ ), the log likelihood function will go to infinity. These singularities provide another example of severe over-fitting of maximum likelihood approach. A further issue is **identifiability** (Casella and Berger, 2002), which is that a  $K$ -component mixture results in  $K!$  equivalent solutions.

### 10.3.3 EM for Gaussian mixtures

An elegant and powerful method for finding the maximum likelihood solutions for models with latent variables is called the **expectation-maximization**, or EM algorithm.

Begin by writing down the conditions that must be satisfied at a maximum of the likelihood function. Setting the derivatives of  $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  10.19 with respect to the means  $\boldsymbol{\mu}_k$  of the Gaussian components to zero, we obtain

$$0 = - \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (10.20)$$

$$= - \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}_k \quad (10.21)$$

$$\therefore N_k = \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (10.22)$$

$$(10.23)$$

where we have made use of  $\frac{de^x}{dx} = e^x$  and 2.52. We can interpret  $N_k$  as the effective number of points assigned to cluster  $k$ .

Now set the derivative of  $\ln(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  with respect to  $\boldsymbol{\Sigma}_k$  to zero:

$$\frac{\partial \ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\Sigma}_k} = \sum_{n=1}^N \frac{\frac{\partial (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))}{\partial \boldsymbol{\Sigma}_k}}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (10.24)$$

$$\therefore \frac{\partial \ln y}{\partial x} = \frac{1}{y} \frac{\partial y}{\partial x} \quad (10.25)$$

$$\therefore \frac{\partial y}{\partial x} = y \frac{\partial \ln y}{\partial x} \quad (10.26)$$

$$\therefore \frac{\partial (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))}{\partial \boldsymbol{\Sigma}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\Sigma}_k} \quad (10.27)$$

$$= \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \frac{\partial \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\Sigma}_k} \quad (10.28)$$

$$\therefore \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\} \quad (10.29)$$

$$\therefore \frac{\partial \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\Sigma}_k} = -\frac{1}{2} \left( \frac{\ln(|\boldsymbol{\Sigma}_k|)}{\partial \boldsymbol{\Sigma}_k} + \frac{\partial (\mathbf{x}_n - \boldsymbol{\mu}_j) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_j)}{\partial \boldsymbol{\Sigma}_k} \right) \quad (10.30)$$

$$\therefore \frac{\partial \det(\mathbf{X})}{\partial \mathbf{X}} = \det(\mathbf{X}) (\mathbf{X}^{-1})^T \quad (10.31)$$

$$\frac{\partial \mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-T} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-T} \quad (10.32)$$

$$\therefore \frac{\ln(|\boldsymbol{\Sigma}_k|)}{\partial \boldsymbol{\Sigma}_k} = \frac{\frac{\partial |\boldsymbol{\Sigma}_k|}{\partial \boldsymbol{\Sigma}_k}}{|\boldsymbol{\Sigma}_k|} = \frac{|\boldsymbol{\Sigma}_k| \boldsymbol{\Sigma}_k^{-T}}{|\boldsymbol{\Sigma}_k|} = \boldsymbol{\Sigma}_k^{-1} \quad (10.33)$$

$$\therefore \frac{\partial \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\Sigma}_k} = -\frac{1}{2} (\boldsymbol{\Sigma}_k^{-1} - \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}) \quad (10.34)$$

$$\therefore \text{derivative} = \sum_{n=1}^N \gamma(z_{nk}) (-1/2) [\boldsymbol{\Sigma}_k^{-1} - \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}] = 0 \quad (10.35)$$

$$\therefore \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \quad (10.36)$$

where we have made use of 2.52 and **The Matrix Cookbook,2012**.

Finally, we maximize  $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  10.19 with respect to the mixing coefficients  $\pi_k$ . Achieve this using a Lagrange multiplier

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right) = 0 \quad (10.37)$$

$$\therefore 0 = \sum_{n=1}^N \frac{\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} + \lambda \quad (10.38)$$

where again we see the appearance of the responsibilities. If we now multiply both sides by  $\pi_k$  and sum over  $k$  making use of constraint  $\sum_{k=1}^K \pi_k = 1$ , we find  $\lambda = -N$ . Then

$$\pi_k = \frac{N_k}{N} \quad (10.39)$$

There do suggest a simple iterative scheme for finding a solution to the maximum likelihood problem.

---

**Algorithm 2:** EM for Gaussian Mixtures

---

1 **repeat**

2     1. Initialize the means  $\boldsymbol{\mu}_k$ , covariances  $\boldsymbol{\Sigma}_k$  and mixing coefficients  $\pi_k$ , and evaluate the initial value of the log likelihood

3     2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \quad (10.40)$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (10.41)$$

4     3. **M step.** Re-estimate the parameters using the current responsibilities

$$\boldsymbol{\mu}_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (10.42)$$

$$\boldsymbol{\Sigma}_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \quad (10.43)$$

$$\pi_k^{new} = \frac{N_k}{N} \quad (10.44)$$

5     4. Evaluate the log likelihood

$$\ln(p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \quad (10.45)$$

6 **until** convergence of either the parameters or the log likelihood

---

## 10.4 Latent Variables View of EM

### 10.4.1 representation

We denote the set of all observed data by  $\mathbf{X}$ , in which the  $n^{th}$  row represents  $x_n^T$ , and denote the set of all latent variables by  $\mathbf{Z}$ , with corresponding rows  $z_n^T$ . The set of all model parameters is denoted by  $\theta$ , so the log likelihood function is given by

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\} \quad (10.46)$$

Replace the sum over  $\mathbf{Z}$  with an integral for continuous latent variables.

$\{\mathbf{X}, \mathbf{Z}\}$  is called the **complete** data set, and  $\{\mathbf{X}\}$  is **incomplete**. We consider complete-data log likelihood functions' expected value under the posterior distribution of the latent variables, corresponding to the E step. In the subsequent M step, maximize this expectation. The expectation of complete-data log likelihood

$$\mathcal{Q}(\theta, \theta^{old}) = \sum_{\mathbf{z}} p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) \quad (10.47)$$

Maximize this

$$\theta^{new} = \arg \max_{\theta} \mathcal{Q}(\theta, \theta^{old}) \quad (10.48)$$

---

#### Algorithm 3: General EM algorithm

---

**Input:**  $p(\mathbf{X}, \mathbf{Z}|\theta), \mathbf{Z}$

**Output:** output

1 **repeat**

2     1. Choose an initial setting for the parameters  $\theta^{old}$ .

3     2. **E step** Evaluate  $p(\mathbf{Z}|\mathbf{X}, \theta^{old})$ .

4     3. **M step** Evaluate  $\theta^{new}$  given by

$$\mathcal{Q}(\theta, \theta^{old}) = \sum_{\mathbf{z}} p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) \quad (10.49)$$

where

$$\theta^{new} = \arg \max_{\theta} \mathcal{Q}(\theta, \theta^{old}) \quad (10.50)$$

5     4.

$$\theta^{old} \leftarrow \theta^{new} \quad (10.51)$$

6 **until converged**

---

When finding MAP(maximum posterior) solutions, maximize the quantity  $\mathcal{Q}(\theta, \theta^{old}) + \ln p(\theta)$  in M step.

## 10.5 The EM Algorithm in General

The **expectation maximization** algorithm, or EM, is a general technique for finding maximum likelihood solutions for probabilistic models having latent variables.

Given the same setting as previous, our goal is to maximize the likelihood function

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \quad (10.52)$$

Replace the summation by integration as appropriate to apply to continuous latent variables.

Introduce a distribution  $q(\mathbf{Z})$  over the latent variables, and decompose the log complete-data likelihood function

$$\therefore \ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) \quad (10.53)$$

$$(10.54)$$

take the expectation over values of  $\mathbf{Z}$  by multiplying both sides by  $q(\mathbf{Z})$  and summing (or integrating) over  $\mathbf{Z}$  we get

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}|\boldsymbol{\theta}) \quad (10.55)$$

$$= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{Z}|\mathbf{X}) \quad (10.56)$$

$$= \sum_{\mathbf{Z}} q(\mathbf{Z}) \left( \ln \frac{p(\mathbf{X}, \mathbf{Z})}{p(\mathbf{Z})} + \ln p(\mathbf{Z}) \right) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{Z}|\mathbf{X}) \quad (10.57)$$

$$= Q + H \quad (10.58)$$

$$= \mathcal{L}(q, \boldsymbol{\theta}) + KL(q \parallel p) \quad (10.59)$$

where

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right\} \quad (10.60)$$

$$KL(q \parallel p) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\} \quad (10.61)$$

Note that  $\mathcal{L}(q, \boldsymbol{\theta})$  is a functional of the distribution  $p(\mathbf{Z})$  and KL is the Kullback-Leibler divergence between  $q(\mathbf{Z})$  and the posterior  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ . KL divergence will be nonnegative according to Jensen's inequality about convex function.

$$\therefore KL(q \parallel p) \geq 0, \text{ with equality if, and only if, } q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) \quad (10.62)$$

$$\therefore \mathcal{L}(q, \boldsymbol{\theta}) \leq \ln p(\mathbf{X}|\boldsymbol{\theta}) \quad (10.63)$$

in other words that  $\mathcal{L}(q, \boldsymbol{\theta})$  is a lower bound on  $\ln p(\mathbf{X}|\boldsymbol{\theta})$ .

The EM is a two-stage **iterative** optimization technique for finding **maximum likelihood** solutions. With the decomposition 10.60, in the E step, the **lower bound**  $\mathcal{L}(q, \boldsymbol{\theta}^{old})$  is maximized with respect to  $q(\mathbf{Z})$  holding  $\boldsymbol{\theta}^{old}$  fixed.

In the subsequent M step, the distribution  $q(\mathbf{Z})$  is held fixed and the lower bound  $\mathcal{L}(q, \boldsymbol{\theta})$  is maximized with respect to  $\boldsymbol{\theta}$ . The nonzero KL divergence causes the log likelihood function increase more than the lower bound, as shown in figure

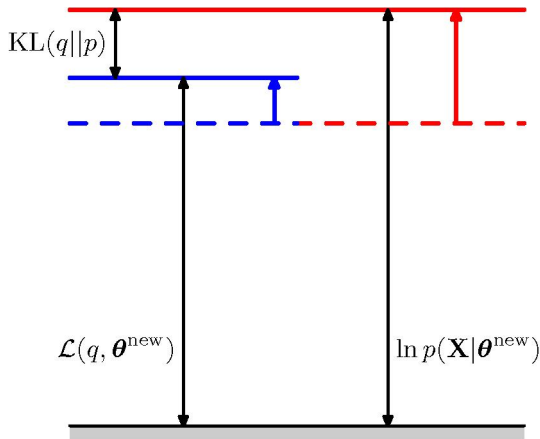


Fig. 10.1: Illustration of the M step of the EM algorithm. The distribution  $q(\mathbf{Z})$  is held fixed and the lower bound  $\mathcal{L}(q, \boldsymbol{\theta})$  is maximized with respect to the parameter vector  $\boldsymbol{\theta}$  to give a revised value  $\boldsymbol{\theta}^{new}$ . Because the KL divergence is nonnegative, this causes the log likelihood  $\ln p(\mathbf{X}|\boldsymbol{\theta})$  to increase by at least as much as the lower bound does

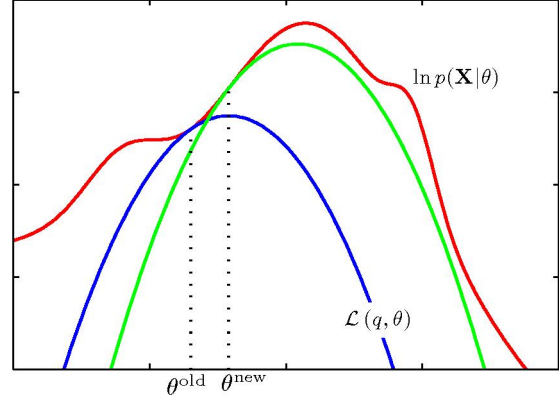
Substitute  $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta^{old})$  into the lower bound, then after the E step,

$$\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \ln p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \quad (10.64)$$

$$= \mathcal{Q}(\theta, \theta^{old}) + const \quad (10.65)$$

where the constant is the negative entropy of the  $q$  distribution.

Fig. 10.2: The EM algorithm involves alternately computing a lower bound on the log likelihood for the current parameter values and then maximizing this bound to obtain the new parameter values. The red curve depicts the incomplete-data log likelihood function to maximize. curve where its value equals the log likelihood at  $\theta^{old}$  indicates the E step. In the M step, the lower bound is maximized giving  $\theta^{new}$ . The shows the subsequent E step, constructing a tangential bound at  $\theta^{new}$



For the particular case of independent, identically distributed data set,  $\mathbf{X} = (x_1^T, x_2^T, \dots, x_n^T)^T$  and  $\mathbf{Z} = (z_1^T, z_2^T, \dots, z_n^T)^T$ . Using the sum and product rules, the posterior in E step take the form

$$p(\mathbf{Z}|\mathbf{X}, \theta) = \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{\sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta)} = \frac{\prod_{n=1}^N p(x_n, z_n|\theta)}{\sum_{\mathbf{Z}} \prod_{n=1}^N p(x_n, z_n|\theta)} = \prod_{n=1}^N p(z_n|x_n, \theta) \quad (10.66)$$

will also factorize with respect to  $n$ .

### 10.5.1 Maximum posterior

Maximize the posterior distribution  $p(\theta|\mathbf{X})$  with EM for models with a prior  $p(\theta)$  over parameters.

$$\therefore p(\theta|\mathbf{X}) = p(\theta, \mathbf{X})/p(\mathbf{X}) \quad (10.67)$$

$$\therefore \ln p(\theta|\mathbf{X}) = \ln p(\theta, \mathbf{X}) - \ln p(\mathbf{X}) = \ln p(\theta) + \ln p(\mathbf{X}|\theta) - \ln p(\mathbf{X}) \quad (10.68)$$

decomposed as

$$\ln p(\theta|\mathbf{X}) = \mathcal{L}(q, \theta) + KL(q \parallel p) + \ln p(\theta) - \ln p(\mathbf{X}) \quad (10.69)$$

$$\geq \mathcal{L}(q, \theta) + \ln p(\theta) - \ln p(\mathbf{X}) \quad (10.70)$$

where  $\ln p(\mathbf{X})$  is constant. We can optimize the right-hand side alternatively with respect to  $q$  and  $\theta$  in E-step and M-step.

### 10.5.2 generalized EM

The **generalized EM**, or GEM, algorithm addresses the problem of an intractable M step. Instead of aiming to maximize  $\mathcal{L}(q, \theta)$  with respect to  $\theta$ , it seeks instead to change the parameters in such a way as to increase its value.

- 1 Nonlinear optimization. conjugate gradients algorithm in M step
- 2 expectation conditional maximization(ECM). Making several constrained optimizations within each M step.

Generalize the E step by performing a **partial**, rather than complete, optimization of  $\mathcal{L}(q, \theta)$  with respect to  $q(\mathbf{Z})$ .





# Chapter 11

## Continuous Latent Variables

### 11.1 Principal Component Analysis

#### 11.1.1 Introduction

Principal Component Analysis is widely used for applications such as dimensionality reduction, lossy data compression, feature extraction, and data visualization. Also known as the Karhunen-Loeve transform. There are two definitions giving rise to the same algorithm. PCA can be defined as the orthogonal projection of the data onto a lower dimensional linear space, known as the principal subspace, such that the variance of the projected data is maximized. Equivalently, it can be defined as the linear projection that minimizes the average projection cost, defined as the linear projection that minimizes the average projection cost, defined as the mean squared distance between the data points and their projections.

#### 11.1.2 Maximum variance formulation

Consider a data set of observations  $\{x_n\}$  where  $n = 1, \dots, N$ , and  $x_n$  is a Euclidean variable with dimensionality  $D$ . Our goal is to project the data onto a space having dimensionality  $M < D$  while maximizing the variance of the projected data. We define the direction of this space using a  $D$ -dimensional unit vector  $\mathbf{u}_1^T \mathbf{u}_1 = 1$ . Each data point  $\mathbf{x}_n$  is then projected onto a scalar value  $\mathbf{u}_1^T \mathbf{x}_n$ . The mean of the projected data is  $\mathbf{u}_1^T \bar{\mathbf{x}}$  where the  $\bar{\mathbf{x}}$  is the sample set mean given by

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad (11.1)$$

and the variance of the projected data is given by

$$\frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}}\}^2 = \frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^T (\mathbf{x}_n - \bar{\mathbf{x}})\}^2 \quad (11.2)$$

$$= \frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^T (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_1\} \quad (11.3)$$

$$= \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 \quad (11.4)$$

where  $\mathbf{S}$  is the data covariance matrix defined by

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \quad (11.5)$$

We now maximize the projected variance  $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$  with respect to  $\mathbf{u}_1$ , which is a constrained maximization to prevent  $\|\mathbf{u}_1\| \rightarrow \infty$ . The appropriate constraint comes from the normalization condition  $\mathbf{u}_1^T \mathbf{u}_1 = 1$ . To enforce this constraint, we introduce a Lagrange multiplier that we shall denote by  $\lambda_1$ , and then make an unconstrained maximization of

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) \quad (11.6)$$

By setting the derivative with respect to  $\mathbf{u}_1$  equal to zero, we see that this quantity will have a stationary point when

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1 \quad (11.7)$$

which says that  $\mathbf{u}_1$  must be an eigenvector of  $\mathbf{S}$ . If we left-multiply by  $\mathbf{u}_1^T$  and make use of  $\mathbf{u}_1^T \mathbf{u}_1 = 1$ , we see that the variance is given by

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1 \quad (11.8)$$

and so the variance will be a maximum when we set  $\mathbf{u}_1$  equal to the eigenvector having the largest eigenvalue  $\lambda_1$ . This eigenvector is known as the first principal component.

### **11.1.3 Minimum-error formulation**

### **11.1.4 Applications of PCA**

### **11.1.5 PCA for high-dimensional data**

## **11.2 Probabilistic PCA**

## **11.3 Kernel PCA**

## **11.4 Nonlinear Latent Variable Models**

## Chapter 12

# Combining Models

### 12.1 Introduction

**committees** Train  $L$  different models and then make predictions using the average of the predictions made by each model.

**boosting** Train multiple models in sequence in which the error function used to train a particular model depends on the performance of the previous models.

**decision trees** Different models are responsible for making predictions in different regions of input space.

**mixtures of experts** Models are viewed as mixture distributions in which the component densities, as well as the mixing coefficients, are conditioned on the input variables.

$$p(t|\mathbf{x}) = \sum_{k=1}^K \pi_k(\mathbf{x}) p(t|\mathbf{x}, k) \quad (12.1)$$

in which  $\pi_k(\mathbf{x}) = p(k|\mathbf{x})$  represents the input-dependent mixing coefficients, and  $k$  indexes the model.

### 12.2 Bayesian Model Averaging

In Bayesian model averaging the whole data set is generated by a single model. By contrast, when we combine multiple models, we see that different data points within the data set can potentially be generated from different values of the latent variable  $z$  and hence by different components.

### 12.3 Committees

The simplest way to construct a committee is to average the predictions of a set of individual models, to cancel the contribution arising from variance and bias.

**Bootstrap** data to introduce variability between the different models within the committee. Suppose we generate  $M$  bootstrap data sets

$$y_{COM}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M y_m(\mathbf{x}) \quad (12.2)$$

where  $m = 1, \dots, M$ . This procedure is known as **bootstrap aggregation** or **bagging**.

### 12.4 Boosting

Here we describe the most widely used form of boosting algorithm: **AdaBoost**, short for 'adaptive boosting'. The base classifiers are known as **weak learners** and are trained in **sequence** using a **weighted form of the data set** in which the weighting coefficient associated with each data point depends on the performance of the previous classifiers.

Consider a two-class classification problem, in which the training data comprises input vectors  $\mathbf{x}_1, \dots, \mathbf{x}_N$  along with corresponding binary target variables  $t_1, \dots, t_N$  where  $t_n \in \{-1, 1\}$ . Each data point is given an associated weighting parameter  $w_n$ , initially set  $1/N$  for all. A base classifier function  $y(\mathbf{x}) \in \{-1, 1\}$

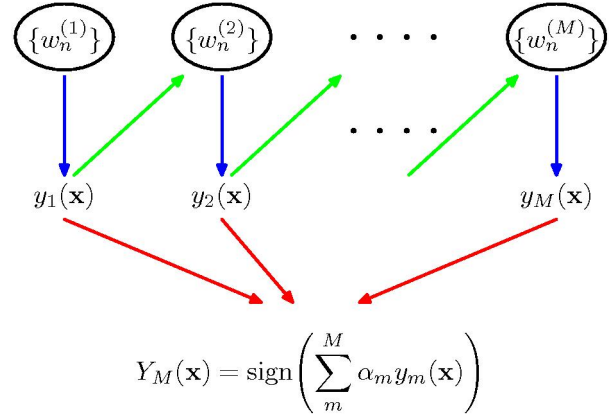


Fig. 12.1: Schematic illustration of boosting framework

**Algorithm 4: AdaBoost****Input:** A set  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}, \{t_1, \dots, t_N\}$ **Output:**  $y(\mathbf{x})$ .

- 1 1. Initialize the data weighting coefficients  $\{w_n\}$  by setting  $w_n^1 = 1/N$  for  $n = 1, \dots, N$ .
- 2 2. **for**  $m = 1, \dots, N$  **do**
- 3     (a) Fit a classifier  $y_m(\mathbf{x})$  to the training data by minimizing the weighted error function

$$J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) \quad (12.3)$$

where  $I(y_m(\mathbf{x}_n) \neq t_n)$  is the indicator function and equals 1 when  $y_m(\mathbf{x}_n) \neq t_n$  and 0 otherwise.

- 4     (b) Evaluate the quantities

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}} \quad (12.4)$$

and then use these to evaluate

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}. \quad (12.5)$$

- 5     (c) Update the data weighting coefficients

$$w_n^{(m+1)} = w_n^{(m)} \exp\{\alpha_m I(y_m(\mathbf{x}_m) \neq t_n)\} \quad (12.6)$$

- 6 3. Make predictions using the final model, which is given by

$$Y_M(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^M \alpha_m y_m(\mathbf{x}) \right) \quad (12.7)$$

**12.4.1 Minimizing exponential error**

Consider the exponential error function defined by

$$E = \sum_{n=1}^N \exp\{-t_n f_m(\mathbf{x}_n)\} \quad (12.8)$$

where  $f_m(\mathbf{x})$  is a classifier defined in terms of a linear combination of base classifiers  $y_l(\mathbf{x})$  of the form

$$f_m(\mathbf{x}) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x}) \quad (12.9)$$

and  $t_n \in \{-1, 1\}$  are the training set target values. Our goal is to minimize  $E$  with respect to the weighting coefficients  $\alpha_l$  and parameters of the base classifiers  $y_l(\mathbf{x})$ .

Separating off the contribution from base classifier  $y_m(\mathbf{x})$ ,

$$E = \sum_{n=1}^N \exp\{-t_n f_m(\mathbf{x}_n)\} \quad (12.10)$$

$$= \sum_{n=1}^N \exp\left\{-t_n \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x})\right\} \quad (12.11)$$

$$= \sum_{n=1}^N \exp\left\{-t_n f_{m-1}(\mathbf{x}_n) - \frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n)\right\} \quad (12.12)$$

$$= \sum_{n=1}^N w_n^{(m)} \exp\left\{-\frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n)\right\} \quad (12.13)$$

where the coefficients  $w_n^{(m)} = \exp\{-t_n f_{m-1}(\mathbf{x}_n)\}$  can be viewed as constants because we are optimizing only  $\alpha_m$  and  $y_m(\mathbf{x})$ . Denote by  $\mathcal{T}_m$  the set of data points correctly classified by  $y_m(\mathbf{x})$  and misclassified points by  $\mathcal{M}_m$ , then we in turn rewrite the error function

$$E = e^{-\alpha_m/2} \sum_{n \in \mathcal{T}_m} w_n^{(m)} + e^{\alpha_m/2} \sum_{n \in \mathcal{M}_m} w_n^{(m)} \quad (12.14)$$

$$= (e^{\alpha_m/2} - e^{-\alpha_m/2}) \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) + e^{-\alpha_m/2} \sum_{n=1}^N w_n^{(m)} \quad (12.15)$$

Then we can minimize this with respect to  $y_m(\mathbf{x}_n)$  and  $\alpha_m$ .

$$\therefore w_n^{(m)} = \exp\{-t_n f_{m-1}(\mathbf{x}_n)\} \quad (12.16)$$

$$\therefore w_n^{(m+1)} = \exp\{-t_n f_m(\mathbf{x}_n)\} \quad (12.17)$$

$$\therefore w_n^{(m+1)} = w_n^{(m)} \exp\left\{-\frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n)\right\} \quad (12.18)$$

Making use of the fact that

$$t_n y_m(\mathbf{x}) = 1 - 2I(y_m(\mathbf{x}_n) \neq t_n) \quad (12.19)$$

we see updates at the next iteration

$$w_n^{(m+1)} = w_n^{(m)} \exp(-\alpha_m/2) \exp\{\alpha_m I(y_m(\mathbf{x}_n) \neq t_n)\} \quad (12.20)$$

Because the term  $\exp(-\alpha_m/2)$  is independent of  $n$ , so can be discarded.

**12.4.2 Error functions for boosting****12.5 Tree-based Models****12.6 Conditional Mixture Models**

# Appendix A

## Matrix Decomposition

### A.1 LU Decomposition

### A.2 QR Decomposition

### A.3 Cholesky Decomposition

### A.4 Eigen value Decomposition

**Definition A.1.** An **eigenvector** of an  $n \times n$  matrix  $A$  is a nonzero vector  $\mathbf{x}$  such that  $A\mathbf{x} = \lambda\mathbf{x}$  for some scalar  $\lambda$ . A scalar  $\lambda$  is called **eigenvalue** of  $A$  if there is a nontrivial solution  $\mathbf{x}$  of  $A\mathbf{x} = \lambda\mathbf{x}$ ; such an  $\mathbf{x}$  is called an eigenvector corresponding to  $\lambda$ <sup>1</sup>.

### A.5 Singular Value Decomposition

#### A.5.1 Definition

**Definition A.2.** Any matrix can be decomposed as follows

$$\underbrace{\mathbf{X}}_{N \times D} = \underbrace{\mathbf{U}}_{N \times N} \underbrace{\mathbf{\Sigma}}_{N \times D} \underbrace{\mathbf{V}^T}_{D \times D} \quad (\text{A.1})$$

where  $\mathbf{U}$  is an  $N \times N$  matrix whose columns are orthonormal (so  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ ),  $\mathbf{V}$  is  $D \times D$  matrix whose rows and columns are orthonormal (so  $\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}_D$ ), and  $\mathbf{\Sigma}$  is a  $N \times D$  matrix containing the  $r = \min(N, D)$  singular values  $\sigma_i \geq 0$  on the main diagonal, with 0s filling the rest of the matrix.

#### A.5.2 Proof

Let  $A$  be an  $m \times n$  matrix. Then  $A^T A$  is symmetric and can be orthogonally diagonalized with eigenvectors. The **singular values** of  $A$  are the square root of the eigenvalues of  $A^T A$ , denoted by  $\sigma_1, \sigma_2, \dots, \sigma_n$ . That is  $\sigma_i = \sqrt{\lambda_i}$  for  $1 \leq i \leq n$ . The eigenvalues are usually arranged so that

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0 \quad (\text{A.2})$$

**Theorem A.1.** Suppose  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$  is an orthogonal basis of  $\mathbb{R}^n$  consisting of eigenvector of  $A^T A$ , arranged so that the corresponding eigenvalues of  $A^T A$  satisfy  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$  and suppose  $A$  has  $r$  nonzero singular values. Then  $A\mathbf{v}_1, \dots, A\mathbf{v}_r$  is an orthogonal basis for  $\text{Col} A$ , and  $\text{rank} A = r$ .

*Proof.* Because  $\mathbf{v}_i$  and  $\mathbf{v}_j$  are orthogonal for  $i \neq j$ ,

$$(A\mathbf{v}_i)^T (A\mathbf{v}_j) = \mathbf{v}_i^T A^T A \mathbf{v}_j = \mathbf{v}_i^T (\lambda_j \mathbf{v}_j) = 0 \quad (\text{A.3})$$

We therefore have

---

<sup>1</sup> An eigenvalue may be zero

$$Av_i = \sigma_i u_i \quad (\text{A.4})$$

For a general vector  $x$ , since eigenvectors are orthogonal unit vectors, we have

$$x = (v_1 \cdot x)v_1 + (v_2 \cdot x)v_2 + \dots + (v_n \cdot x)v_n \quad (\text{A.5})$$

This means that

$$Mx = (v_1 \cdot x)Mv_1 + (v_2 \cdot x)Mv_2 + \dots + (v_n \cdot x)Mv_n \quad (\text{A.6})$$

$$Mx = (v_1 \cdot x)\sigma_1 u_1 + (v_2 \cdot x)\sigma_2 u_2 + \dots + (v_n \cdot x)\sigma_n u_n \quad (\text{A.7})$$

Remember that dot product can be computed using the vector transpose

$$v \cdot u = v^T u \quad (\text{A.8})$$

which leads to

$$Mx = u_1 \sigma_1 v_1^T x + u_2 \sigma_2 v_2^T x + \dots + u_n \sigma_n v_n^T x \quad (\text{A.9})$$

$$M = u_1 \sigma_1 v_1^T + u_2 \sigma_2 v_2^T + \dots + u_n \sigma_n v_n^T \quad (\text{A.10})$$

And this is usually expressed by writing

$$M = U \Sigma V^T \quad (\text{A.11})$$

As for  $u_i$ , we have

$$\begin{aligned} & \begin{cases} (A^T A)v_i = \lambda_i v_i \\ Av_i = \sigma_i u_i \end{cases} \\ & \Rightarrow A^T \sigma_i u_i = \lambda_i v_i \\ & \Rightarrow \sigma_i A^T u_i = \lambda_i v_i \\ & \Rightarrow (AA^T)u_i = \sigma_i A v_i = \lambda_i u_i \end{aligned}$$

So we can see that  $u_i$  is the eigenvector of symmetric matrix  $AA^T$ , and  $v_i$  is the eigenvector of symmetric matrix  $A^T A$ . In summary,  $u_i$  and  $v_i$  are the **left-eigenvector** and **right-eigenvectors** of matrix  $A$ .

## A.5.3 Application

### A.5.3.1 Principal Component Analysis

The projection vectors for principal component projection are the left-eigenvectors



## Appendix B

### Optimization methods

#### B.1 Convexity

**Definition B.1. (Convex set)** We say a set  $\mathcal{S}$  is convex if for any  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{S}$ , we have

$$\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \mathcal{S}, \forall \lambda \in [0, 1] \quad (\text{B.1})$$

**Definition B.2. (Convex function)** A function  $f(\mathbf{x})$  is called convex if its **epigraph** (the set of points above the function) defines a convex set. Equivalently, a function  $f(\mathbf{x})$  is called convex if it is defined on a convex set and if, for any  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{S}$ , and any  $\lambda \in [0, 1]$ , we have

$$f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) \leq \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2) \quad (\text{B.2})$$

**Definition B.3.** A function  $f(\mathbf{x})$  is said to be **strictly convex** if the inequality is strict

$$f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) < \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2) \quad (\text{B.3})$$

**Definition B.4.** A function  $f(\mathbf{x})$  is said to be (strictly) **concave** if  $-f(\mathbf{x})$  is (strictly) convex.

**Theorem B.1.** If  $f(x)$  is twice differentiable on  $[a, b]$  and  $f''(x) \geq 0$  on  $[a, b]$  then  $f(x)$  is convex on  $[a, b]$ .

**Proposition B.1.**  $\log(x)$  is strictly convex on  $(0, \infty)$ .

Intuitively, a (strictly) convex function has a bowl shape, and hence has a unique global minimum  $\mathbf{x}^*$  corresponding to the bottom of the bowl. Hence its second derivative must be positive everywhere,  $\frac{d^2}{dx^2} f(x) > 0$ . A twice-continuously differentiable, multivariate function  $f$  is convex iff its Hessian is positive definite for all  $\mathbf{x}$ . In the machine learning context, the function  $f$  often corresponds to the NLL.

Models where the NLL is convex are desirable, since this means we can always find the globally optimal MLE. We will see many examples of this later in the book. However, many models of interest will not have concave likelihoods. In such cases, we will discuss ways to derive locally optimal parameter estimates.

#### B.2 Gradient descent

##### B.2.1 Stochastic gradient descent

---

**Algorithm 5:** Stochastic gradient descent

---

**input** : Training data  $\mathcal{D} = \{(\mathbf{x}_i, y_i) | i = 1 : N\}$   
**output**: A linear model:  $y_i = \boldsymbol{\theta}^T \mathbf{x}$

```

1  $\mathbf{w} \leftarrow 0; b \leftarrow 0; k \leftarrow 0;$ 
2 while no mistakes made within the for loop do
3   for  $i \leftarrow 1$  to  $N$  do
4     if  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \leq 0$  then
5        $\mathbf{w} \leftarrow \mathbf{w} + \eta y_i \mathbf{x}_i;$ 
6        $b \leftarrow b + \eta y_i;$ 
7        $k \leftarrow k + 1;$ 
```

---

## B.2.2 Batch gradient descent

## B.2.3 Line search

The **line search**<sup>2</sup> approach first finds a descent direction along which the objective function  $f$  will be reduced and then computes a step size that determines how far  $\mathbf{x}$  should move along that direction. The descent direction can be computed by various methods, such as gradient descent(Section B.2), Newton's method(Section B.4) and Quasi-Newton method(Section B.5). The step size can be determined either exactly or inexactly.

## B.2.4 Momentum term

## B.3 Lagrange duality

### B.3.1 Primal form

Consider the following, which we'll call the **primal** optimization problem:

$$xyz \tag{B.4}$$

### B.3.2 Dual form

## B.4 Newton's method

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \mathbf{g}_k^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^T \mathbf{H}_k (\mathbf{x} - \mathbf{x}_k)$$

$$\text{where } \mathbf{g}_k \triangleq \mathbf{g}(\mathbf{x}_k) = f'(\mathbf{x}_k), \mathbf{H}_k \triangleq \mathbf{H}(\mathbf{x}_k),$$

$$\mathbf{H}(\mathbf{x}) \triangleq \left[ \frac{\partial^2 f}{\partial x_i \partial x_j} \right]_{D \times D} \quad (\text{Hessian matrix})$$

$$f'(\mathbf{x}) = \mathbf{g}_k + \mathbf{H}_k (\mathbf{x} - \mathbf{x}_k) = 0 \Rightarrow \tag{B.5}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}_k^{-1} \mathbf{g}_k \tag{B.6}$$

---

**Algorithm 6:** Newtons method for minimizing a strictly convex function

---

```

1 Initialize  $\mathbf{x}_0$ 
2 while (!convergency) do
3   Evaluate  $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$ 
4   Evaluate  $\mathbf{H}_k = \nabla^2 f(\mathbf{x}_k)$ 
5    $\mathbf{d}_k = -\mathbf{H}_k^{-1} \mathbf{g}_k$ 
6   Use line search to find step size  $\eta_k$  along  $\mathbf{d}_k$ 
7    $\mathbf{x}_{k+1} = \mathbf{x}_k + \eta_k \mathbf{d}_k$ 

```

---

<sup>2</sup> [http://en.wikipedia.org/wiki/Line\\_search](http://en.wikipedia.org/wiki/Line_search)

## B.5 Quasi-Newton method

From Equation B.5 we can infer out the **quasi-Newton condition** as follows:

$$\begin{aligned}
 f'(x) - g_k &= H_k(x - x_k) \\
 g_{k-1} - g_k &= H_k(x_{k-1} - x_k) \Rightarrow \\
 g_k - g_{k-1} &= H_k(x_k - x_{k-1}) \\
 g_{k+1} - g_k &= H_{k+1}(x_{k+1} - x_k) \quad (\text{quasi-Newton condition})
 \end{aligned} \tag{B.7}$$

The idea is to replace  $H_k^{-1}$  with a approximation  $B_k$ , which satisfies the following properties:

1.  $B_k$  must be symmetric
2.  $B_k$  must satisfies the quasi-Newton condition, i.e.,  $g_{k+1} - g_k = B_{k+1}(x_{k+1} - x_k)$ .

Let  $y_k = g_{k+1} - g_k$ ,  $\delta_k = x_{k+1} - x_k$ , then

$$B_{k+1}y_k = \delta_k \tag{B.8}$$

3. Subject to the above,  $B_k$  should be as close as possible to  $B_{k-1}$ .

Note that we did not require that  $B_k$  be positive definite. That is because we can show that it must be positive definite if  $B_{k-1}$  is. Therefore, as long as the initial Hessian approximation  $B_0$  is positive definite, all  $B_k$  are, by induction.

### B.5.1 DFP

Updating rule:

$$B_{k+1} = B_k + P_k + Q_k \tag{B.9}$$

From Equation B.8 we can get

$$B_{k+1}y_k = B_k y_k + P_k y_k + Q_k y_k = \delta_k$$

To make the equation above establish, just let

$$\begin{aligned}
 P_k y_k &= \delta_k \\
 Q_k y_k &= -B_k y_k
 \end{aligned}$$

In DFP algorithm,  $P_k$  and  $Q_k$  are

$$P_k = \frac{\delta_k \delta_k^T}{\delta_k^T y_k} \tag{B.10}$$

$$Q_k = -\frac{B_k y_k y_k^T B_k}{y_k^T B_k y_k} \tag{B.11}$$

### B.5.2 BFGS

Use  $B_k$  as a approximation to  $H_k$ , then the quasi-Newton condition becomes

$$B_{k+1} \delta_k = y_k \tag{B.12}$$

The updating rule is similar to DFP, but  $P_k$  and  $Q_k$  are different. Let

$$\begin{aligned} P_k \delta_k &= y_k \\ Q_k \delta_k &= -B_k \delta_k \end{aligned}$$

Then

$$P_k = \frac{y_k y_k^T}{y_k^T \delta_k} \quad (\text{B.13})$$

$$Q_k = -\frac{B_k \delta_k \delta_k^T B_k}{\delta_k^T B_k \delta_k} \quad (\text{B.14})$$

### B.5.3 Broyden

Broyden's algorithm is a linear combination of DFP and BFGS.

# Glossary

Use the template *glossary.tex* together with the Springer document class SVMono (monograph-type books) or SVMult (edited books) to style your glossary in the Springer layout.

**glossary term** Write here the description of the glossary term. Write here the description of the glossary term. Write here the description of the glossary term.

**glossary term** Write here the description of the glossary term. Write here the description of the glossary term. Write here the description of the glossary term.

**glossary term** Write here the description of the glossary term. Write here the description of the glossary term. Write here the description of the glossary term.

**glossary term** Write here the description of the glossary term. Write here the description of the glossary term. Write here the description of the glossary term.

**glossary term** Write here the description of the glossary term. Write here the description of the glossary term. Write here the description of the glossary term.