

<https://github.com/soulmachine/machine-learning-cheat-sheet>  
soulmachine@gmail.com

# Machine Learning Cheat Sheet

Classical equations, diagrams and tricks in machine learning

December 9, 2013

©2013 soulmachine

Except where otherwise noted, This document is licensed under a Creative Commons Attribution-ShareAlike 3.0 Unported (CC BY-SA3.0) license

(<http://creativecommons.org/licenses/by/3.0/>).

# Preface

This cheat sheet contains many classical equations and diagrams on machine learning, which will help you quickly recall knowledge and ideas on machine learning.

The cheat sheet will also appeal to someone who is preparing for a job interview related to machine learning.

This cheat sheet has three significant advantages:

1. Strong typed. Compared to programming languages, mathematical formulas are weakly typed. For example,  $X$  can be a set, a random variable, or a matrix. This causes difficulty in understanding the meaning of formulas. In this cheat sheet, I try my best to standardize symbols used, see section §.
2. More parentheses. In machine learning, authors are prone to omit parentheses, brackets and braces, this usually causes ambiguity in mathematical formulas. In this cheat sheet, I use parentheses(brackets and braces) at where they are needed, to make formulas easy to understand.
3. Less thinking jumps. In many books, authors are prone to omit some steps that are trivial in his option. But it often makes readers get lost in the middle way of derivation.

At Tsinghua University, May 2013

*soulmachine*



# Contents

<b>Contents</b>	v
<b>Notation</b>	xiii
<b>1 Introduction</b>	1
1.1 Types of machine learning	1
1.2 Three elements of a machine learning model	1
1.2.1 Representation	1
1.2.2 Evaluation	1
1.2.3 Optimization	3
1.3 Cross validation	3
<b>2 Probability</b>	5
2.1 Frequentists vs. Bayesians	5
2.2 A brief review of probability theory	5
2.2.1 Basic concepts	5
2.2.2 Multivariate random variables	6
2.2.3 Bayes rule	7
2.2.4 Independence and conditional independence	7
2.2.5 Quantiles	7
2.2.6 Mean and variance	7
2.3 Some common discrete distributions	8
2.3.1 The Bernoulli and binomial distributions	8
2.3.2 The multinoulli and multinomial distributions	8
2.3.3 The Poisson distribution	9
2.3.4 The empirical distribution	9
2.4 Some common continuous distributions	9
2.4.1 Gaussian (normal) distribution	9
2.4.2 Student's t-distribution	10
2.4.3 The Laplace distribution	11
2.4.4 The gamma distribution	11
2.4.5 The beta distribution	12
2.4.6 Pareto distribution	12
2.5 Joint probability distributions	13
2.5.1 Covariance and correlation	14
2.5.2 Multivariate Gaussian distribution	15
2.5.3 Multivariate Student's t-distribution	16
2.5.4 Dirichlet distribution	16
2.6 Transformations of random variables	16

2.6.1	Linear transformations	18
2.6.2	General transformations	18
2.6.3	Central limit theorem	18
2.7	Monte Carlo approximation	19
2.8	Information theory	19
2.8.1	Entropy	19
2.8.2	KL divergence	19
2.8.3	Mutual information	20
<b>3</b>	<b>Generative models for discrete data</b>	<b>23</b>
3.1	Generative classifier	23
3.2	Bayesian concept learning	23
3.2.1	Likelihood	23
3.2.2	Prior	23
3.2.3	Posterior	24
3.2.4	Posterior predictive distribution	25
3.3	The beta-binomial model	26
3.3.1	Likelihood	26
3.3.2	Prior	26
3.3.3	Posterior	27
3.3.4	Posterior predictive distribution	28
3.4	The Dirichlet-multinomial model	29
3.4.1	Likelihood	29
3.4.2	Prior	30
3.4.3	Posterior	30
3.4.4	Posterior predictive distribution	30
3.5	Naive Bayes classifiers	31
3.5.1	Optimization	31
3.5.2	Using the model for prediction	32
3.5.3	The log-sum-exp trick	32
3.5.4	Feature selection using mutual information	33
3.5.5	Classifying documents using bag of words	33
<b>4</b>	<b>Gaussian Models</b>	<b>35</b>
4.1	Basics	35
4.1.1	MLE for a MVN	36
4.1.2	Maximum entropy derivation of the Gaussian *	36
4.2	Gaussian discriminant analysis	36
4.2.1	Quadratic discriminant analysis (QDA)	37
4.2.2	Linear discriminant analysis (LDA)	37
4.2.3	Two-class LDA	39
4.2.4	MLE for discriminant analysis	40
4.2.5	Strategies for preventing overfitting	40
4.2.6	Regularized LDA *	41
4.2.7	Diagonal LDA	41
4.2.8	Nearest shrunken centroids classifier *	41
4.3	Inference in jointly Gaussian distributions	41
4.3.1	Statement of the result	41
4.3.2	Examples	42
4.4	Inferring the parameters of an MVN	42
4.4.1	Posterior distribution of $\mu$	42
4.4.2	Posterior distribution of $\Sigma$ *	42
4.4.3	Posterior distribution of $\mu$ and $\Sigma$ *	42
4.4.4	Sensor fusion with unknown precisions *	42

<b>5</b>	<b>Bayesian statistics</b>	43
5.1	Introduction	43
5.2	Summarizing posterior distributions	43
5.2.1	MAP estimation	43
5.2.2	Credible intervals	45
5.2.3	Inference for a difference in proportions	45
5.3	Bayesian model selection	46
5.3.1	Bayesian Occam's razor	46
5.3.2	Computing the marginal likelihood (evidence)	47
5.3.3	Bayes factors	50
5.4	Priors	50
5.4.1	Uninformative priors	50
5.4.2	Robust priors	50
5.4.3	Mixtures of conjugate priors	50
5.5	Hierarchical Bayes	51
5.6	Empirical Bayes	51
5.7	Bayesian decision theory	51
5.7.1	Bayes estimators for common loss functions	51
5.7.2	The false positive vs false negative tradeoff	53
<b>6</b>	<b>Frequentist statistics</b>	55
6.1	Sampling distribution of an estimator	55
6.1.1	Bootstrap	55
6.1.2	Large sample theory for the MLE *	55
6.2	Frequentist decision theory	55
6.3	Desirable properties of estimators	56
6.4	Empirical risk minimization	56
6.4.1	Regularized risk minimization	56
6.4.2	Structural risk minimization	56
6.4.3	Estimating the risk using cross validation	56
6.4.4	Upper bounding the risk using statistical learning theory *	56
6.4.5	Surrogate loss functions	56
6.5	Pathologies of frequentist statistics *	56
<b>7</b>	<b>Linear Regression</b>	57
7.1	Introduction	57
7.2	Representation	57
7.3	MLE	57
7.3.1	OLS	58
7.3.2	SGD	59
7.4	Ridge regression(MAP)	60
7.4.1	Basic idea	60
7.4.2	Numerically stable computation *	61
7.4.3	Connection with PCA *	61
7.4.4	Regularization effects of big data	61
7.5	Bayesian linear regression	61
<b>8</b>	<b>Logistic Regression</b>	63
8.1	Representation	63
8.2	Optimization	63
8.2.1	MLE	63
8.2.2	MAP	64
8.3	Multinomial logistic regression	64
8.3.1	Representation	64

8.3.2	MLE	64
8.3.3	MAP	65
8.4	Bayesian logistic regression	66
8.4.1	Laplace approximation	66
8.4.2	Derivation of the BIC	66
8.4.3	Gaussian approximation for logistic regression	66
8.4.4	Approximating the posterior predictive	66
8.4.5	Residual analysis (outlier detection) *	66
8.5	Online learning and stochastic optimization	66
8.6	Generative vs discriminative classifiers	66
8.6.1	Pros and cons of each approach	66
8.6.2	Dealing with missing data	67
8.6.3	Fishers linear discriminant analysis (FLDA) *	69
<b>9</b>	<b>Generalized linear models and the exponential family</b>	<b>71</b>
9.1	The exponential family	71
9.1.1	Definition	71
9.1.2	Examples	72
9.1.3	Log partition function	73
9.1.4	MLE for the exponential family	74
9.1.5	Bayes for the exponential family	75
9.1.6	Maximum entropy derivation of the exponential family *	75
9.2	Generalized linear models (GLMs)	75
9.2.1	Basics	75
9.3	Probit regression	75
9.4	Multi-task learning	75
<b>10</b>	<b>Directed graphical models (Bayes nets)</b>	<b>77</b>
10.1	Introduction	77
10.1.1	Chain rule	77
10.1.2	Conditional independence	77
10.1.3	Graphical models	77
10.1.4	Directed graphical model	77
10.2	Examples	78
10.2.1	Naive Bayes classifiers	78
10.2.2	Markov and hidden Markov models	79
10.3	Inference	79
10.4	Learning	79
10.4.1	Learning from complete data	79
10.4.2	Learning with missing and/or latent variables	80
10.5	Conditional independence properties of DGMs	80
10.5.1	d-separation and the Bayes Ball algorithm (global Markov properties)	80
10.5.2	Other Markov properties of DGMs	80
10.5.3	Markov blanket and full conditionals	80
10.5.4	Multinoulli Learning	80
10.6	Influence (decision) diagrams *	81
<b>11</b>	<b>Mixture models and the EM algorithm</b>	<b>83</b>
11.1	Latent variable models	83
11.2	Mixture models	83
11.2.1	Mixtures of Gaussians	84
11.2.2	Mixtures of multinoullis	84
11.2.3	Using mixture models for clustering	85
11.2.4	Mixtures of experts	85



11.3	Parameter estimation for mixture models	87
11.3.1	Unidentifiability	87
11.3.2	Computing a MAP estimate is non-convex	87
11.4	The EM algorithm	87
11.4.1	Introduction	87
11.4.2	Basic idea	87
11.4.3	EM for GMMs	88
11.4.4	EM for mixture of experts	91
11.4.5	EM for DGMs with hidden variables	91
11.4.6	EM for the Student distribution *	91
11.4.7	EM for probit regression *	91
11.4.8	Derivation of the $Q$ function	91
11.4.9	Convergence of the EM Algorithm *	92
11.4.10	Generalization of EM Algorithm *	93
11.4.11	Online EM	94
11.4.12	Other EM variants *	94
11.5	Model selection for latent variable models	94
11.5.1	Model selection for probabilistic models	95
11.5.2	Model selection for non-probabilistic methods	95
11.6	Fitting models with missing data	95
11.6.1	EM for the MLE of an MVN with missing data	96
12	Latent linear models	97
13	Sparse linear models	99
14	Kernels	101
14.1	Introduction	101
14.2	Kernel functions	101
14.2.1	RBF kernels	101
14.2.2	TF-IDF kernels	102
14.2.3	Mercer (positive definite) kernels	102
14.2.4	Linear kernels	103
14.2.5	Matern kernels	103
14.2.6	String kernels	103
14.2.7	Pyramid match kernels	104
14.2.8	Kernels derived from probabilistic generative models	104
14.3	Using kernels inside GLMs	105
14.3.1	Kernel machines	105
14.3.2	L1VMs, RVMs, and other sparse vector machines	106
14.4	The kernel trick	107
14.4.1	Kernelized KNN	107
14.4.2	Kernelized K-medoids clustering	108
14.4.3	Kernelized ridge regression	108
14.4.4	Kernel PCA	109
14.5	Support vector machines (SVMs)	109
14.5.1	SVMs for classification	109
14.5.2	SVMs for regression	111
14.5.3	Choosing $C$	112
14.5.4	A probabilistic interpretation of SVMs	113
14.5.5	Summary of key points	113
14.6	Comparison of discriminative kernel methods	113
14.7	Kernels for building generative models	114

<b>15</b>	<b>Perceptron</b>	115
15.1	Representation	115
15.2	Evaluation	115
15.3	Optimization	116
15.3.1	Primal form	116
15.3.2	Dual form	116
<b>16</b>	<b>K-Nearest Neighbors</b>	119
16.1	Representation	119
16.2	Evaluation	119
16.3	Optimization	119
<b>17</b>	<b>K-Means Clustering</b>	121
17.1	Representation	121
17.2	Evaluation	121
17.3	Optimization	121
17.4	Tricks	122
17.4.1	Choosing $k$	122
17.4.2	Choosing the initial centroids(seeds)	122
17.5	Reference	122
<b>18</b>	<b>Decision Tree</b>	123
<b>19</b>	<b>AdaBoost</b>	125
19.1	Representation	125
19.2	Evaluation	125
19.3	Optimization	125
19.3.1	Input	125
19.3.2	Output	125
19.3.3	Algorithm	126
19.4	The upper bound of the training error of AdaBoost	126
<b>20</b>	<b>Hidden markov Model</b>	127
<b>21</b>	<b>Conditional Random Field</b>	129
<b>A</b>	<b>Optimization methods</b>	131
A.1	Convexity	131
A.2	Gradient descent	132
A.2.1	Stochastic gradient descent	132
A.2.2	Batch gradient descent	132
A.2.3	Line search	132
A.2.4	Momentum term	132
A.3	Lagrange duality	132
A.3.1	Primal form	132
A.3.2	Dual form	132
A.4	Newton's method	132
A.5	Quasi-Newton method	133
A.5.1	DFP	133
A.5.2	BFGS	134
A.5.3	Broyden	134
	<b>Glossary</b>	135

## List of Contributors

Wei Zhang

PhD candidate at the Institute of Software, Chinese Academy of Sciences (ISCAS), Beijing, P.R.CHINA, e-mail: [zh3feng@gmail.com](mailto:zh3feng@gmail.com), has written chapters of Naive Bayes and SVM.

Fei Pan

Master at Beijing University of Technology, Beijing, P.R.CHINA, e-mail: [example@gmail.com](mailto:example@gmail.com), has written chapters of KMeans, AdaBoost.

Yong Li

PhD candidate at the Institute of Automation of the Chinese Academy of Sciences (CASIA), Beijing, P.R.CHINA, e-mail: [liyong3forever@gmail.com](mailto:liyong3forever@gmail.com), has written chapters of Logistic Regression.

Jiankou Li

PhD candidate at the Institute of Software, Chinese Academy of Sciences (ISCAS), Beijing, P.R.CHINA, e-mail: [lijiankoucoco@163.com](mailto:lijiankoucoco@163.com), has written chapters of BayesNet.



# 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
$ 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
$\text{tr}(\mathbf{X})$	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
$\text{diag}(\mathbf{x})$	Diagonal matrix made from vector $\mathbf{x}$
$\text{diag}(\mathbf{X})$	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
$cov[X]$	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 $\Lambda = \Sigma^{-1}$
$mode[X]$	Most probable value of $X$
$\mu$	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 $\ell_1$ regularizer
$\phi(x)$	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 $i$

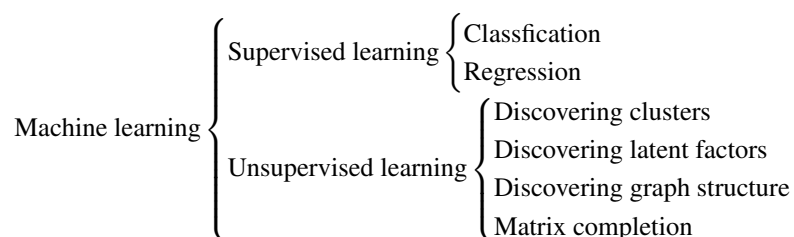
---



# Chapter 1

## Introduction

### 1.1 Types of machine learning



### 1.2 Three elements of a machine learning model

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

#### 1.2.1 Representation

In supervised learning, a model must be represented as a conditional probability distribution  $P(y|\mathbf{x})$  (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.

#### 1.2.2 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.

##### 1.2.2.1 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  $L(y_i, \hat{y})$ .

---

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

The following is some common loss functions:

1. 0-1 loss function  $L(Y, f(X)) = I(Y, f(X)) = \begin{cases} 1, & Y \neq f(X) \\ 0, & Y = f(X) \end{cases}$
2. Quadratic loss function  $L(Y, f(X)) = (Y - f(X))^2$
3. Absolute loss function  $L(Y, f(X)) = |Y - f(X)|$
4. Logarithmic loss function  $L(Y, P(Y|X)) = -\log P(Y|X)$

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

$$R_{\text{exp}}(f) = E_p[L(Y, f(X))] = \int_{X \times Y} L(y, f(x)) P(x, y) dx dy \quad (1.1)$$

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

$$R_{\text{emp}}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) \quad (1.2)$$

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 §15).

### 1.2.2.2 ERM and SRM

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

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

**Definition 1.5.** Structural risk

$$R_{\text{smp}}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f) \quad (1.4)$$

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

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

---

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

### 1.2.3 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.

## 1.3 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.

---

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



## Chapter 2

# Probability

### 2.1 Frequentists vs. Bayesians

what is 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 ought 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 A brief review of probability theory

#### 2.2.1 Basic concepts

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

**Discrete random variable:**  $X$  can take on any value from a finite or countably infinite set.

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

##### 2.2.1.1 CDF

$$F(x) \triangleq P(X \leq x) = \begin{cases} \sum_{u \leq x} p(u), & \text{discrete random variable} \\ \int_{-\infty}^x f(u) du, & \text{continuous random variable} \end{cases} \quad (2.1)$$

##### 2.2.1.2 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 Multivariate random variables

### 2.2.2.1 Joint CDF

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

$$F(x, y) \triangleq P(X \leq x, Y \leq y) = \begin{cases} \sum_{u \leq x, v \leq y} p(u, v), & \text{descrete} \\ \int_{-\infty}^x \int_{-\infty}^y f(u, v) du dv, & \text{continuous} \end{cases} \quad (2.2)$$

**product rule:**

$$p(X, Y) = P(X|Y)P(Y) \quad (2.3)$$

**Chain rule:**

$$p(X_{1:N}) = p(X_1)p(X_2|X_1)p(X_3|X_2, X_1) \dots p(X_N|X_{1:N-1}) \quad (2.4)$$

### 2.2.2.2 Marginal distribution

**Marginal CDF:**

$$\begin{cases} F_X(x) \triangleq F(x, +\infty) = \begin{cases} \sum_{x_i \leq x} P(X = x_i) = \sum_{x_i \leq x} \sum_{j=1}^{+\infty} P(X = x_i, Y = y_j), & \text{descrete} \\ \int_{-\infty}^x f_X(u) du = \int_{-\infty}^x \int_{-\infty}^{+\infty} f(u, v) du dv, & \text{continuous} \end{cases} \\ F_Y(y) \triangleq F(+\infty, y) = \begin{cases} \sum_{y_j \leq y} p(Y = y_j) = \sum_{i=1}^{+\infty} \sum_{y_j \leq y} P(X = x_i, Y = y_j), & \text{descrete} \\ \int_{-\infty}^y f_Y(v) dv = \int_{-\infty}^{+\infty} \int_{-\infty}^y f(u, v) du dv, & \text{continuous} \end{cases} \end{cases} \quad (2.5)$$

**Marginal PMF or PDF:**

$$\begin{cases} \begin{cases} P(X = x_i) = \sum_{j=1}^{+\infty} P(X = x_i, Y = y_j), & \text{descrete} \\ f_X(x) = \int_{-\infty}^{+\infty} f(x, y) dy, & \text{continuous} \end{cases} \\ \begin{cases} p(Y = y_j) = \sum_{i=1}^{+\infty} P(X = x_i, Y = y_j), & \text{descrete} \\ f_Y(y) = \int_{-\infty}^{+\infty} f(x, y) dx, & \text{continuous} \end{cases} \end{cases} \quad (2.6)$$

### 2.2.2.3 Conditional distribution

**Conditional PMF:**

$$p(X = x_i | Y = y_j) = \frac{p(X = x_i, Y = y_j)}{p(Y = y_j)} \text{ if } p(Y) > 0 \quad (2.7)$$

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

**Conditional PDF:**

$$f_{X|Y}(x|y) = \frac{f(x, y)}{f_Y(y)} \quad (2.8)$$

<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 Bayes rule

$$p(Y = y|X = x) = \frac{p(X = x, Y = y)}{p(X = x)} = \frac{p(X = x|Y = y)p(Y = y)}{\sum_{y'} p(X = x|Y = y')p(Y = y')} \quad (2.9)$$

### 2.2.3.1 Example: Generative classifiers

$$p(y = c|\mathbf{x}, \boldsymbol{\theta}) = \frac{p(\mathbf{x}|y = c, \boldsymbol{\theta})p(y = c|\boldsymbol{\theta})}{\sum_{c'} p(\mathbf{x}|y = c', \boldsymbol{\theta})p(y = c'|\boldsymbol{\theta})} \quad (2.10)$$

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**.

## 2.2.4 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.,

$$X \perp Y = P(X, Y) = P(X)P(Y) \quad (2.11)$$

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

$$X \perp Y|Z = P(X, Y|Z) = P(X|Z)P(Y|Z) \quad (2.12)$$

## 2.2.5 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**.

## 2.2.6 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:

$$\text{var}[X] = \mathbb{E}[(X - \mu)^2] = \int (x - \mu)^2 p(x)dx \quad (2.13)$$

$$= \int x^2 p(x)dx - 2\mu \int xp(x)dx + \mu^2 \int p(x)dx = \mathbb{E}[X^2] - \mu^2 \quad (2.14)$$

from which we derive the useful result

$$\mathbb{E}[X^2] = \sigma^2 + \mu^2 \quad (2.15)$$

The **standard deviation** is defined as

$$\text{std}[X] \triangleq \sqrt{\text{var}[X]} \quad (2.16)$$

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

## 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

$$\text{Ber}(x|\theta) \triangleq \theta^{\mathbb{I}(x=1)}(1-\theta)^{\mathbb{I}(x=0)} \quad (2.17)$$

**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

$$\text{Bin}(k|n, \theta) \triangleq \binom{n}{k} \theta^k (1-\theta)^{n-k} \quad (2.18)$$

### 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:

$$p(\mathbf{x}) \triangleq \prod_{k=1}^K \theta_k^{\mathbb{I}(x_k=1)} \quad (2.19)$$

**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

$$p(\mathbf{x}) \triangleq \binom{n}{x_1 \dots x_K} \prod_{k=1}^K \theta_k^{x_k} \quad \text{where} \quad \binom{n}{x_1 \dots x_K} \triangleq \frac{n!}{x_1! x_2! \dots x_K!} \quad (2.20)$$

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 2.2 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  $X \sim \text{Poi}(\lambda)$ , if its pmf is

$$p(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!} \quad (2.21)$$

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	$X \sim \text{Cat}(\boldsymbol{\theta})$	$\mathbf{x} \in \{0, 1\}^K, \sum_{k=1}^K x_k = 1$	$\prod_{k=1}^K \theta_j^{\mathbb{I}(x_j=1)}$	-	-
Multinomial	$X \sim \text{Mu}(n, \boldsymbol{\theta})$	$\mathbf{x} \in \{0, 1, \dots, n\}^K, \sum_{k=1}^K x_k = n$	$\binom{n}{x_1 \dots x_K} \prod_{k=1}^K \theta_j^{x_j}$	-	-
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

$$F_n(x) \triangleq \frac{1}{N} \sum_{i=1}^N \mathbb{I}(x_i \leq x) \quad (2.22)$$

## 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.

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

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

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.
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.

## 2.4.2 Student's t-distribution

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

Name	Written as	$f(x)$	$\mathbb{E}[X]$	mode	$\text{var}[X]$
Student's t-distribution $X \sim \mathcal{T}(\mu, \sigma^2, \nu)$	$\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})}$	$\left[1 + \frac{1}{\nu} \left(\frac{x-\mu}{\sigma}\right)^2\right]^{-\frac{\nu+1}{2}}$	$\mu$	$\mu$	$\frac{\nu\sigma^2}{\nu-2}$

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

$$\Gamma(x) \triangleq \int_0^\infty u^{x-1} e^{-u} du \quad (2.23)$$

$\mu$  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.

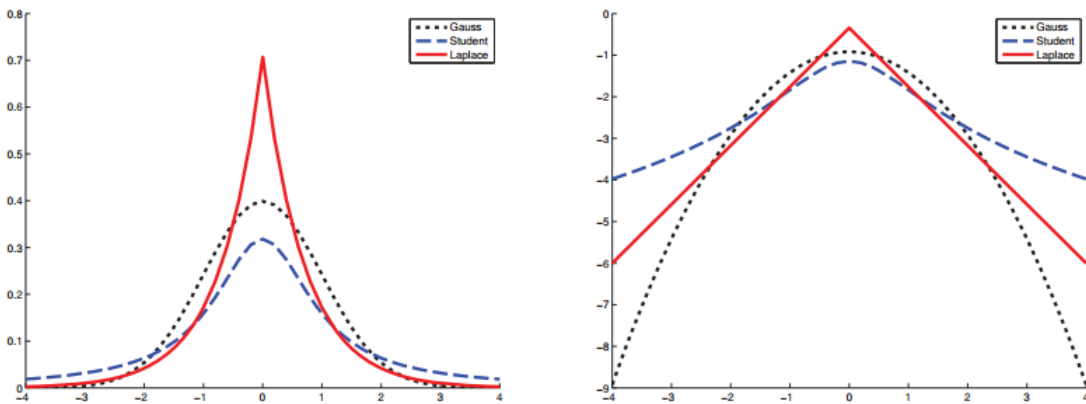


Fig. 2.1: (a) The pdfs for a  $\mathcal{N}(0, 1)$ ,  $\mathcal{T}(0, 1, 1)$  and  $\text{Lap}(0, 1/\sqrt{2})$ . The mean is 0 and the variance is 1 for both the Gaussian and Laplace. The mean and variance of the Student is undefined when  $\nu = 1$ . (b) Log of these pdfs. Note that the Student distribution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

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).

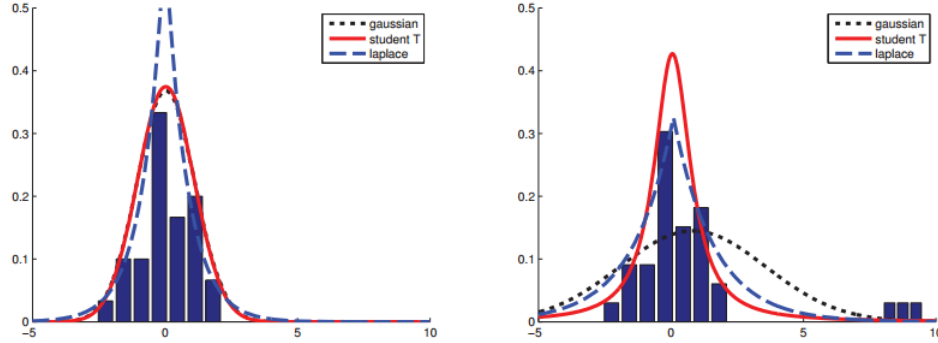


Fig. 2.2: Illustration of the effect of outliers on fitting Gaussian, Student and Laplace distributions. (a) No outliers (the Gaussian and Student curves are on top of each other). (b) With outliers. We see that the Gaussian is more affected by outliers than the Student and Laplace distributions.

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.

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

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.

### 2.4.4 The gamma distribution

Table 2.6: Summary of gamma distribution

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

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

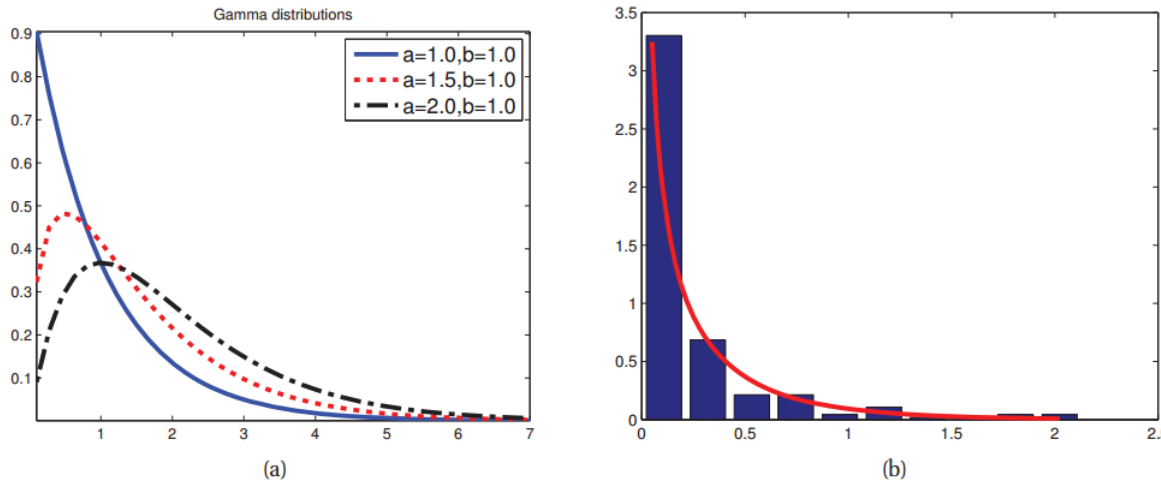


Fig. 2.3: Some  $\text{Ga}(a, b = 1)$  distributions. If  $a \leq 1$ , the mode is at 0, otherwise it is  $> 0$ . As we increase the rate  $b$ , we reduce the horizontal scale, thus squeezing everything leftwards and upwards. (b) An empirical pdf of some rainfall data, with a fitted Gamma distribution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

## 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,

$$B(a, b) \triangleq \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \quad (2.24)$$

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.

## 2.4.6 Pareto distribution

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**).

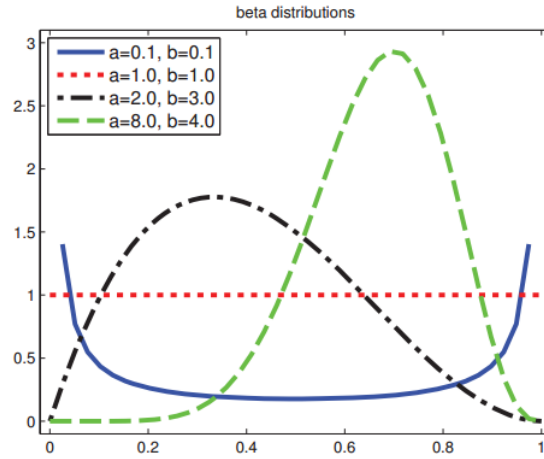
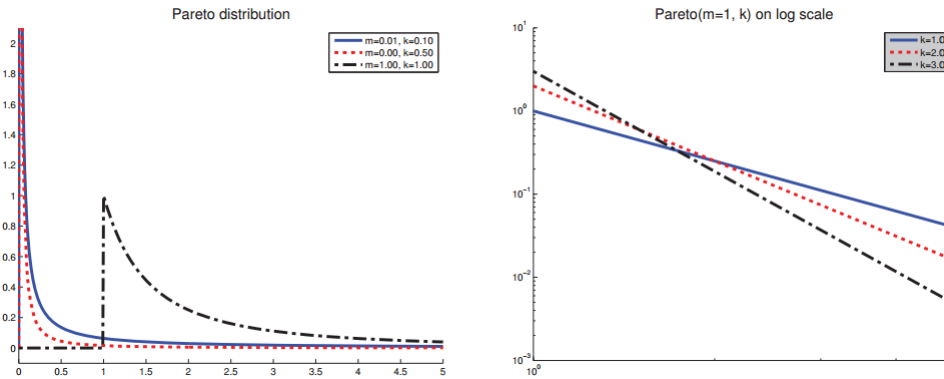


Fig. 2.4: Some beta distributions.

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 $k > 2$

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 Joint probability distributions

Given a **multivariate random variable** or **random vector**<sup>7</sup>  $X \in \mathbb{R}^D$ , the **joint probability distribution**<sup>8</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).

<sup>7</sup> [http://en.wikipedia.org/wiki/Multivariate\\_random\\_variable](http://en.wikipedia.org/wiki/Multivariate_random_variable)

<sup>8</sup> [http://en.wikipedia.org/wiki/Joint\\_probability\\_distribution](http://en.wikipedia.org/wiki/Joint_probability_distribution)

### 2.5.1 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

$$\text{Cov}[X, Y] \triangleq \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \quad (2.25)$$

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

$$\text{Cov}[X] \triangleq \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T] \quad (2.26)$$

$$= \begin{pmatrix} \text{var}[X_1] & \text{Cov}[X_1, X_2] & \cdots & \text{Cov}[X_1, X_D] \\ \text{Cov}[X_2, X_1] & \text{var}[X_2] & \cdots & \text{Cov}[X_2, X_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[X_D, X_1] & \text{Cov}[X_D, X_2] & \cdots & \text{var}[X_D] \end{pmatrix} \quad (2.27)$$

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

$$\text{corr}[X, Y] \triangleq \frac{\text{Cov}[X, Y]}{\sqrt{\text{var}[X] \text{var}[Y]}} \quad (2.28)$$

A **correlation matrix** has the form

$$\mathbf{R} \triangleq \begin{pmatrix} \text{corr}[X_1, X_1] & \text{corr}[X_1, X_2] & \cdots & \text{corr}[X_1, X_D] \\ \text{corr}[X_2, X_1] & \text{corr}[X_2, X_2] & \cdots & \text{corr}[X_2, X_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{corr}[X_D, X_1] & \text{corr}[X_D, X_2] & \cdots & \text{corr}[X_D, X_D] \end{pmatrix} \quad (2.29)$$

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

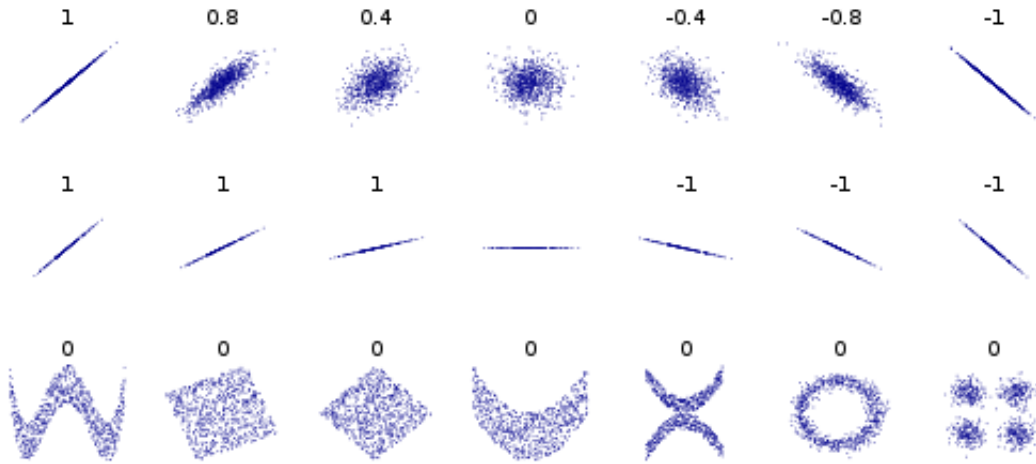


Fig. 2.6: Several sets of  $(x, y)$  points, with the Pearson correlation coefficient of  $x$  and  $y$  for each set. Note that the correlation reflects the noisiness and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many aspects of nonlinear relationships (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.

## 2.5.2 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:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (2.30)$$

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  $\boldsymbol{\Sigma}$  is symmetric). A diagonal covariance matrix has  $D$  parameters, and has 0s in the off-diagonal terms. A spherical or isotropic covariance,  $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_D$ , has one free parameter.

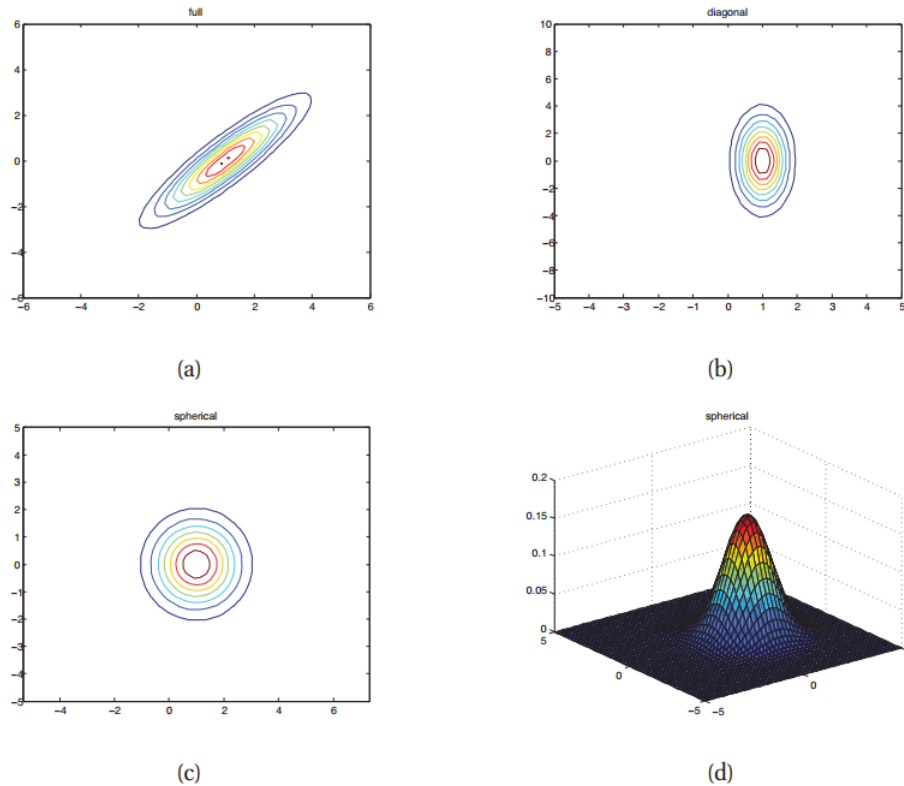


Fig. 2.7: We show the level sets for 2d Gaussians. (a) A full covariance matrix has elliptical contours. (b) A diagonal covariance matrix is an axis aligned ellipse. (c) A spherical covariance matrix has a circular shape. (d) Surface plot for the spherical Gaussian in (c).

### 2.5.3 Multivariate Student's t-distribution

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

$$\mathcal{T}(x|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu) \triangleq \frac{\Gamma(\frac{\nu+D}{2})}{\Gamma(\frac{\nu}{2})} \frac{|\boldsymbol{\Sigma}|^{-\frac{1}{2}}}{(\nu\pi)^{\frac{D}{2}}} \left[ 1 + \frac{1}{\nu} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]^{-\frac{\nu+D}{2}} \quad (2.31)$$

$$= \frac{\Gamma(\frac{\nu+D}{2})}{\Gamma(\frac{\nu}{2})} \frac{|\boldsymbol{\Sigma}|^{-\frac{1}{2}}}{(\nu\pi)^{\frac{D}{2}}} \left[ 1 + (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{V}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]^{-\frac{\nu+D}{2}} \quad (2.32)$$

where  $\boldsymbol{\Sigma}$  is called the scale matrix (since it is not exactly the covariance matrix) and  $\mathbf{V} = \nu\boldsymbol{\Sigma}$ . 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

$$\text{mean} = \boldsymbol{\mu}, \text{mode} = \boldsymbol{\mu}, \text{Cov} = \frac{\nu}{\nu-2} \boldsymbol{\Sigma} \quad (2.33)$$

### 2.5.4 Dirichlet distribution

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

$$S_K = \left\{ \mathbf{x} : 0 \leq x_k \leq 1, \sum_{k=1}^K x_k = 1 \right\} \quad (2.34)$$

The pdf is defined as follows:

$$\text{Dir}(\mathbf{x}|\boldsymbol{\alpha}) \triangleq \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K x_k^{\alpha_k-1} \mathbb{I}(\mathbf{x} \in S_K) \quad (2.35)$$

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

$$B(\boldsymbol{\alpha}) \triangleq \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\alpha_0)} \text{ where } \alpha_0 \triangleq \sum_{k=1}^K \alpha_k \quad (2.36)$$

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

$$\mathbb{E}(x_k) = \frac{\alpha_k}{\alpha_0}, \text{mode}[x_k] = \frac{\alpha_k - 1}{\alpha_0 - K}, \text{var}[x_k] = \frac{\alpha_k(\alpha_0 - \alpha_k)}{\alpha_0^2(\alpha_0 + 1)} \quad (2.37)$$

## 2.6 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.



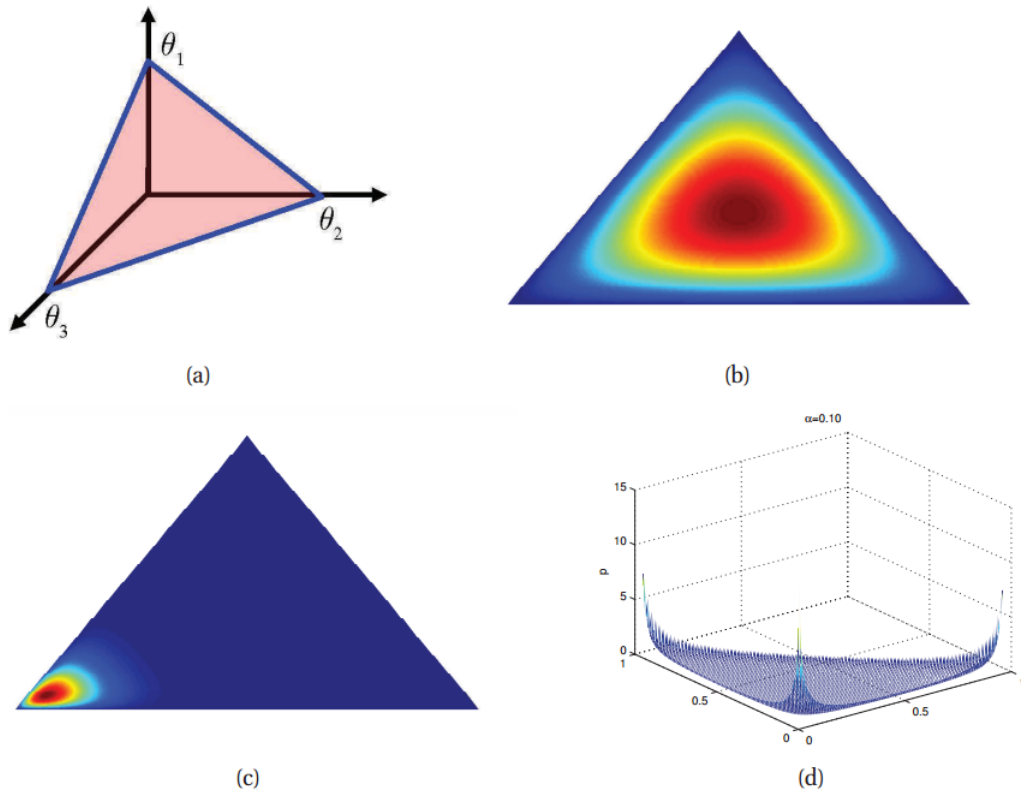


Fig. 2.8: (a) The Dirichlet distribution when  $K = 3$  defines a distribution over the simplex, which can be represented by the triangular surface. Points on this surface satisfy  $0 \leq \theta_k \leq 1$  and  $\sum_{k=1}^K \theta_k = 1$ . (b) Plot of the Dirichlet density when  $\alpha = (2, 2, 2)$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

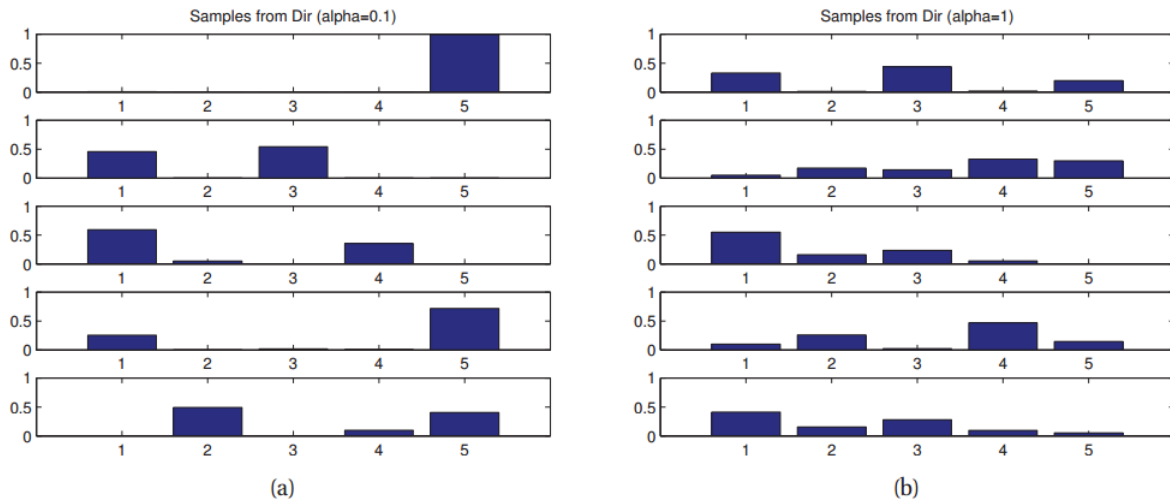


Fig. 2.9: Samples from a 5-dimensional symmetric Dirichlet distribution for different parameter values. (a)  $\alpha = (0.1, \dots, 0.1)$ . This results in very sparse distributions, with many 0s. (b)  $\alpha = (1, \dots, 1)$ . This results in more uniform (and dense) distributions.

### 2.6.1 Linear transformations

Suppose  $g(\cdot)$  is a linear function:

$$g(\mathbf{x}) = A\mathbf{x} + b \quad (2.38)$$

First, for the mean, we have

$$\mathbb{E}[\mathbf{y}] = \mathbb{E}[A\mathbf{x} + b] = A\mathbb{E}[\mathbf{x}] + b \quad (2.39)$$

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

For the covariance, we have

$$\text{Cov}[\mathbf{y}] = \text{Cov}[A\mathbf{x} + b] = A\Sigma A^T \quad (2.40)$$

### 2.6.2 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  $f(x) = y$ :

$$p_Y(y) = \sum_{x:g(x)=y} p_X(x) \quad (2.41)$$

If  $X$  is continuous, we cannot use Equation (2.41) since  $p_X(x)$  is a density, not a pmf, and we cannot sum up densities. Instead, we work with cdfs, and write

$$F_Y(y) = P(Y \leq y) = P(g(X) \leq y) = \int_{g(X) \leq y} f_X(x) dx \quad (2.42)$$

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

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right| \quad (2.43)$$

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  $p_Y(y) = \frac{1}{2}y^{-\frac{1}{2}}$ .

### 2.6.3 Central limit theorem

**Theorem 2.1.** Given  $N$  random variables  $X_1, X_2, \dots, X_N$ , each variable is *independent and identically distributed*<sup>9</sup> (iid for short), and each has the same mean  $\mu$  and variance  $\sigma^2$ , then

$$\frac{\sum_{i=1}^n X_i - N\mu}{\sqrt{N}\sigma} \sim \mathcal{N}(0, 1) \quad (2.44)$$

this can also be written as

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{N}} \sim \mathcal{N}(0, 1) \text{ where } \bar{X} \triangleq \frac{1}{N} \sum_{i=1}^n X_i \quad (2.45)$$

<sup>9</sup> [http://en.wikipedia.org/wiki/Independent\\_identically\\_distributed](http://en.wikipedia.org/wiki/Independent_identically_distributed)

## 2.7 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>10</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:

$$\mathbb{E}[g(X)] = \int g(x)p(x)dx \approx \frac{1}{S} \sum_{s=1}^S f(x_s), \text{ where } x_s \sim p(X) \quad (2.46)$$

This is called **Monte Carlo integration**<sup>11</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.

## 2.8 Information theory

### 2.8.1 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

$$\mathbb{H}(X) \triangleq - \sum_{k=1}^K p(X=k) \log_2 p(X=k) \quad (2.47)$$

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.

### 2.8.2 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:

$$\mathbb{KL}(P||Q) \triangleq \sum_x p(x) \log_2 \frac{p(x)}{q(x)} \quad (2.48)$$

where the sum gets replaced by an integral for pdfs<sup>12</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

<sup>10</sup> [http://en.wikipedia.org/wiki/Monte\\_Carlo\\_method](http://en.wikipedia.org/wiki/Monte_Carlo_method)

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

<sup>12</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  $q = 0.5p_1 + 0.5p_2$

as zero because  $\lim_{x \rightarrow 0} x \ln x$ . We can rewrite this as

$$\mathbb{KL}(p||q) \triangleq \sum_x p(x) \log_2 p(x) - \sum_{k=1}^K p(x) \log_2 q(x) = \mathbb{H}(p) - \mathbb{H}(p, q) \quad (2.49)$$

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

$$\mathbb{H}(p, q) = \sum_x p(x) \log_2 q(x) \quad (2.50)$$

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 §2.8.1, 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 2.2. (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*.

## 2.8.3 Mutual information

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

$$\mathbb{I}(X; Y) \triangleq \mathbb{KL}(P(X, Y) || P(X)P(Y)) = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \quad (2.51)$$

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:

$$\mathbb{I}(X; Y) = \mathbb{H}(X) - \mathbb{H}(X|Y) \quad (2.52)$$

$$= \mathbb{H}(Y) - \mathbb{H}(Y|X) \quad (2.53)$$

$$= \mathbb{H}(X) + \mathbb{H}(Y) - \mathbb{H}(X, Y) \quad (2.54)$$

$$= \mathbb{H}(X, Y) - \mathbb{H}(X|Y) - \mathbb{H}(Y|X) \quad (2.55)$$

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>13</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  $X$ .

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

$$PMI(x, y) \triangleq \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)} \quad (2.56)$$

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:

<sup>13</sup> [http://en.wikipedia.org/wiki/Mutual\\_information](http://en.wikipedia.org/wiki/Mutual_information)

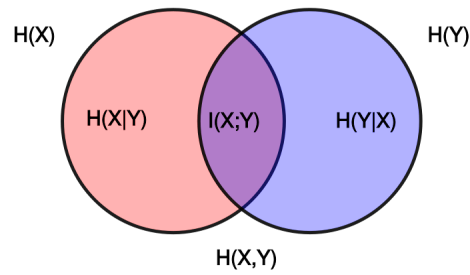


Fig. 2.10: Individual  $\mathbb{H}(X)$ ,  $\mathbb{H}(Y)$ , joint  $\mathbb{H}(X, Y)$ , and conditional entropies for a pair of correlated subsystems  $X, Y$  with mutual information  $\mathbb{I}(X; Y)$ .

$$PMI(x, y) = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)} \quad (2.57)$$

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

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

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  $p(y = c)$ .

#### 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  $C$ .

For pedagogical purposes, we will consider a very simple example of concept learning called the **number game**, based on part of Josh Tenenbaums PhD thesis (Tenenbaum 1999), see Figure 3.1.

##### 3.2.1 Likelihood

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

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>14</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.

---

<sup>14</sup> [http://en.wikipedia.org/wiki/Occam%27s\\_razor](http://en.wikipedia.org/wiki/Occam%27s_razor)

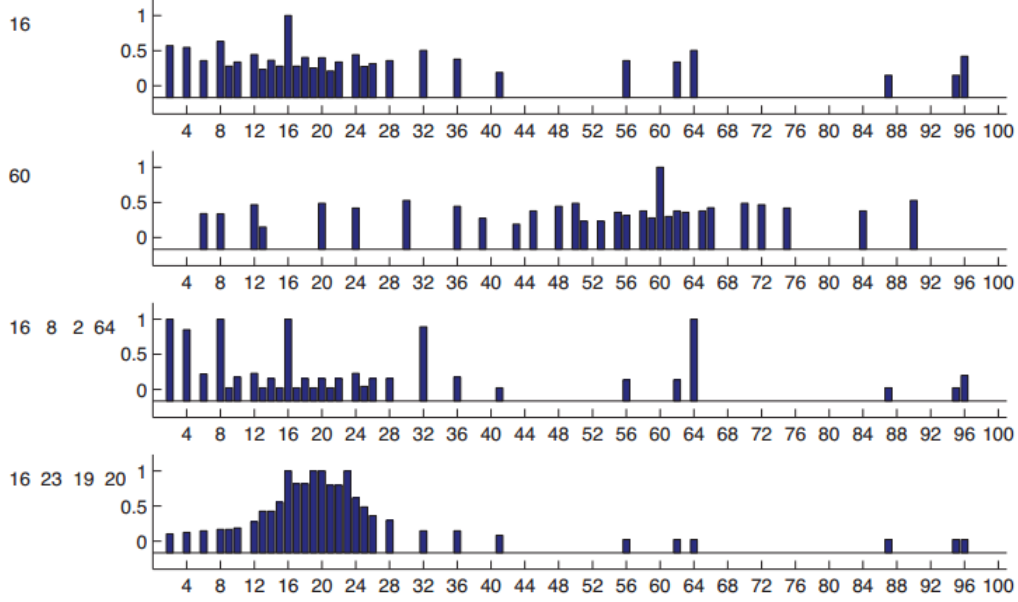


Fig. 3.1: Empirical predictive distribution averaged over 8 humans in the number game. First two rows: after seeing  $\mathcal{D} = 16$  and  $\mathcal{D} = 60$ . This illustrates diffuse similarity. Third row: after seeing  $\mathcal{D} = 16, 8, 2, 64$ . This illustrates rule-like behaviour (powers of 2). Bottom row: after seeing  $\mathcal{D} = 16, 23, 19, 20$ . This illustrates focussed similarity (numbers near 20). Source: Figure 5.5 of (Tenenbaum 1999).

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.

So, what prior should we use? For illustration purposes, let us use a simple prior which puts uniform probability on 30 simple arithmetical concepts, such as even numbers, odd numbers, prime numbers, numbers ending in 9, etc. We also include two unnatural concepts, namely powers of 2, plus 37 and powers of 2, except 32, but give them low prior weight. See Figure 3.2(a) for a plot of this prior.

### 3.2.3 Posterior

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

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

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  $h$ .

Figure 3.2 plots the prior, likelihood and posterior after seeing  $\mathcal{D} = \{16\}$ . We see that the posterior is a combination of prior and likelihood. In the case of most of the concepts, the prior is uniform, so the posterior is proportional to the likelihood. However, the unnatural concepts of powers of 2, plus 37 and powers of 2, except 32 have low posterior support, despite having high likelihood, due to the low prior. Conversely, the concept of odd numbers has low posterior support, despite having a high prior, due to the low likelihood.

Figure 3.3 plots the prior, likelihood and posterior after seeing  $\mathcal{D} = \{16, 8, 2, 64\}$ . Now the likelihood is much more peaked on the powers of two concept, so this dominates the posterior.

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,



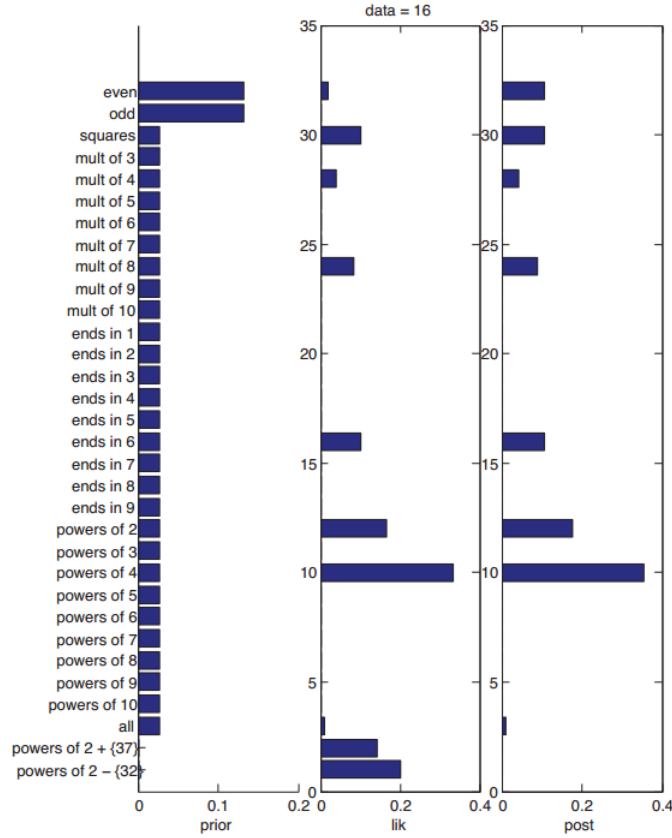


Fig. 3.2: Prior, likelihood and posterior for  $\mathcal{D} = \{16\}$ . Based on (Tenenbaum 1999).

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

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>15</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.

$$p(\tilde{x}|\mathcal{D}) \triangleq \mathbb{E}_{h|\mathcal{D}}[p(\tilde{x}|h)] = \begin{cases} \sum_h p(\tilde{x}|h)p(h|\mathcal{D}) & , \text{ discrete parameters} \\ \int_h p(\tilde{x}|h)p(h|\mathcal{D})dh & , \text{ continuous parameters} \end{cases} \quad (3.7)$$

This is just a weighted average of the predictions of each individual hypothesis and is called **Bayes model averaging**(Hoeting et al. 1999). This is illustrated in Figure 3.4. The dots at the bottom show the predictions from each

<sup>15</sup> [http://en.wikipedia.org/wiki/Posterior\\_predictive\\_distribution](http://en.wikipedia.org/wiki/Posterior_predictive_distribution)

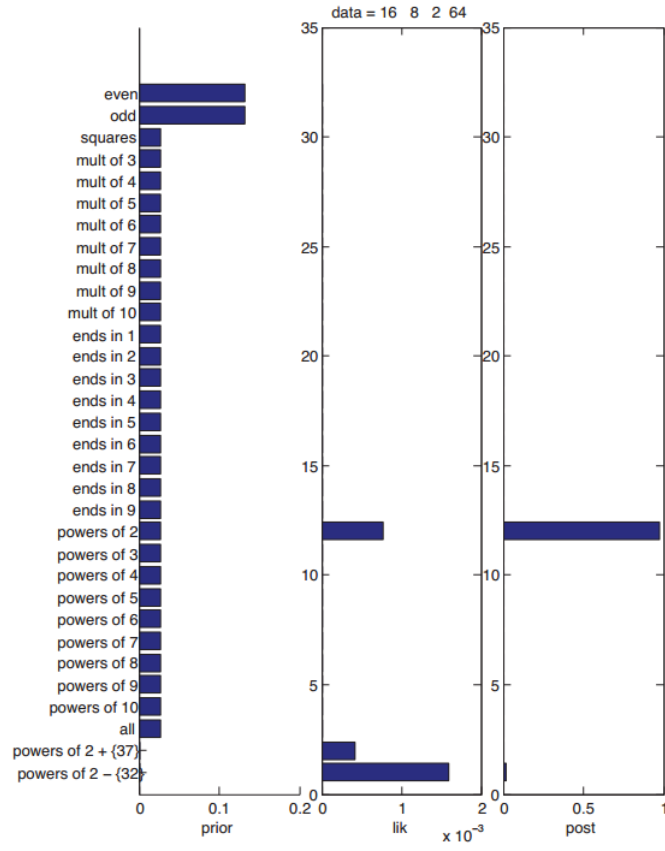


Fig. 3.3: Prior, likelihood and posterior for  $\mathcal{D} = \{16, 8, 2, 64\}$ . Based on (Tenenbaum 1999).

hypothesis; the vertical curve on the right shows the weight associated with each hypothesis. If we multiply each row by its weight and add up, we get the distribution at the top.

### 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**.

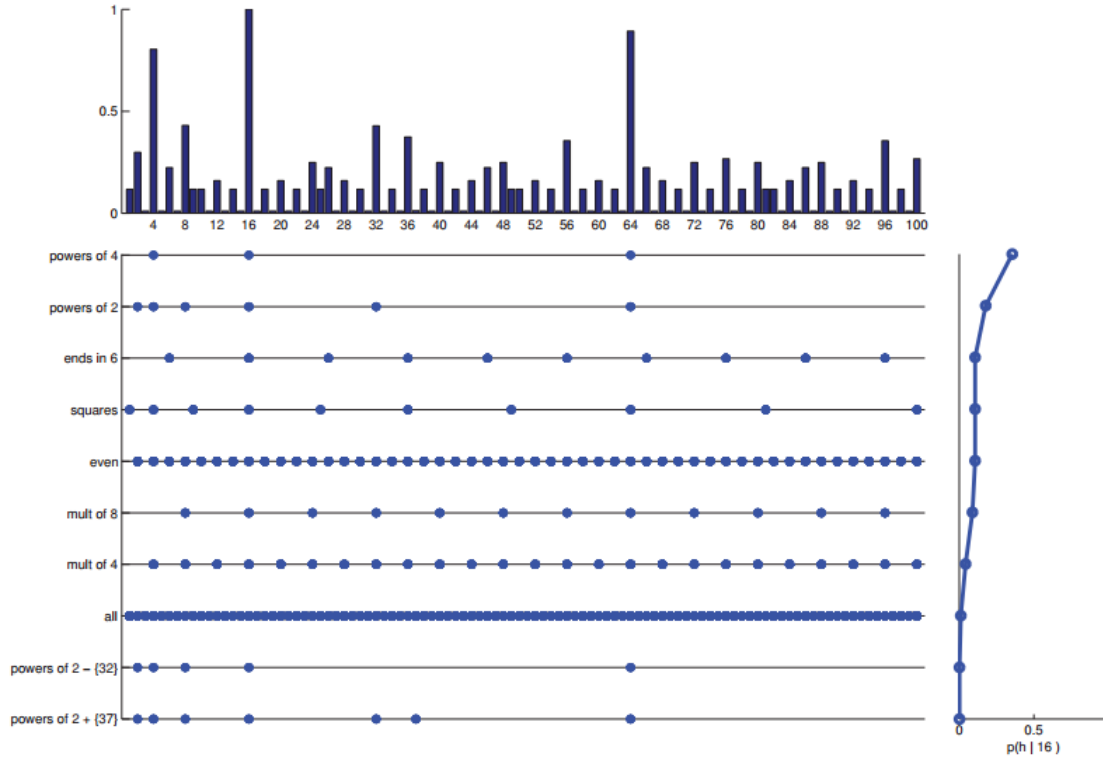


Fig. 3.4: Posterior over hypotheses and the corresponding predictive distribution after seeing one example,  $\mathcal{D} = \{16\}$ . A dot means this number is consistent with this hypothesis. The graph  $p(h|\mathcal{D})$  on the right is the weight given to hypothesis  $h$ . By taking a weighed sum of dots, we get  $p(\tilde{x}|\mathcal{D})$ (top). Based on Figure 2.9 of (Tenenbaum 1999).

### 3.3.3 Posterior

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

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

$$\begin{aligned} p(\theta|\mathcal{D}_a, \mathcal{D}_b) &= p(\theta, \mathcal{D}_b|\mathcal{D}_a)p(\mathcal{D}_a) \\ &\propto p(\theta, \mathcal{D}_b|\mathcal{D}_a) \\ &= p(\mathcal{D}_b, \theta|\mathcal{D}_a) \\ &= p(\mathcal{D}_b|\theta)p(\theta|\mathcal{D}_a) \end{aligned}$$

Combine Equation 3.10 and 2.18

$$\begin{aligned} &= \text{Bin}(N_1^b|\theta, N_1^a + N_0^b) \text{Beta}(\theta|N_1^a + a, N_0^a + b) \\ &= \text{Beta}(\theta|N_1^a + N_1^b + a, N_0^a + N_0^b + b) \end{aligned}$$

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

$$\text{var}(\theta|\mathcal{D}) \approx \frac{N_1 N_0}{N^2} = \frac{\hat{\theta}_{MLE}(1 - \hat{\theta}_{MLE})}{N} \quad (3.15)$$

## 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

$$p(\tilde{x}|\mathcal{D}) = \int_0^1 p(\tilde{x}|\theta)p(\theta|\mathcal{D})d\theta \quad (3.16)$$

$$= \int_0^1 \theta \text{Beta}(\theta|a, b)d\theta = \mathbb{E}[\theta|\mathcal{D}] = \frac{a}{a + b} \quad (3.17)$$

### 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**

$$p(\tilde{x}|\mathcal{D}) = \frac{N_1 + 1}{N_0 + N_1 + 1} \quad (3.18)$$

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

$$p(\tilde{x}|\mathcal{D}) = \int_0^1 \text{Bin}(\tilde{x}|M, \theta) \text{Beta}(\theta|a, b) d\theta \quad (3.19)$$

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

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

$$\int_0^1 \theta^{\tilde{x}} (1 - \theta)^{M - \tilde{x}} \theta^{a-1} (1 - \theta)^{b-1} d\theta = B(\tilde{x} + a, M - \tilde{x} + b) \quad (3.21)$$

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

$$Bb(x|a, b, M) \triangleq \binom{M}{x} \frac{B(x + a, M - x + b)}{B(a, b)} \quad (3.22)$$

This distribution has the following mean and variance

$$\text{mean} = M \frac{a}{a + b}, \text{var} = \frac{Mab}{(a + b)^2} \frac{a + b + M}{a + b + 1} \quad (3.23)$$

This process is illustrated in Figure 3.5. 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 3.5(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.

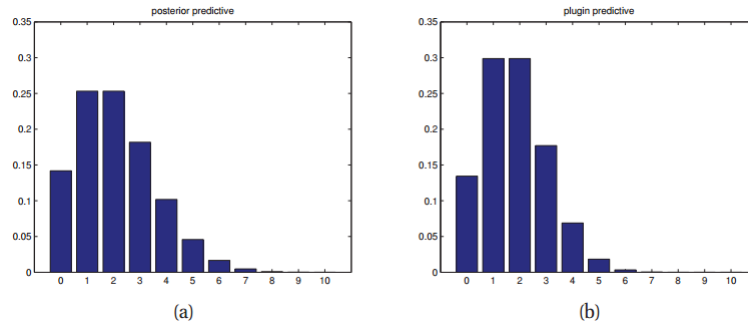


Fig. 3.5: (a) Posterior predictive distributions after seeing  $N_1 = 3, N_0 = 17$ . (b) MAP estimation.

## 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

$$p(\mathcal{D}|\boldsymbol{\theta}) = \binom{N}{N_1 \dots N_K} \prod_{k=1}^K \theta_k^{N_k} \quad \text{where } N_k = \sum_{i=1}^N \mathbb{I}(y_i = k) \quad (3.24)$$

almost the same as Equation (2.20).

### 3.4.2 Prior

$$\text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \mathbb{I}(\boldsymbol{\theta} \in S_K) \quad (3.25)$$

### 3.4.3 Posterior

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (3.26)$$

$$\propto \prod_{k=1}^K \theta_k^{N_k} \theta_k^{\alpha_k - 1} = \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1} \quad (3.27)$$

$$= \text{Dir}(\boldsymbol{\theta}|\alpha_1 + N_1, \dots, \alpha_K + N_K) \quad (3.28)$$

From Equation (2.37), the MAP estimate is given by

$$\hat{\theta}_k = \frac{N_k + \alpha_k - 1}{N + \alpha_0 - K} \quad (3.29)$$

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

$$\hat{\theta}_k = \frac{N_k}{N} \quad (3.30)$$

### 3.4.4 Posterior predictive distribution

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

$$p(X = j|\mathcal{D}) = \int p(X = j|\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta} \quad (3.31)$$

$$= \int p(X = j|\boldsymbol{\theta}_j) \left[ \int p(\boldsymbol{\theta}_{-j}, \boldsymbol{\theta}_j|\mathcal{D})d\boldsymbol{\theta}_{-j} \right] d\boldsymbol{\theta}_j \quad (3.32)$$

$$= \int \boldsymbol{\theta}_j p(\boldsymbol{\theta}_j|\mathcal{D})d\boldsymbol{\theta}_j = \mathbb{E}[\boldsymbol{\theta}_j|\mathcal{D}] = \frac{\boldsymbol{\alpha}_j + N_j}{\boldsymbol{\alpha}_0 + N} \quad (3.33)$$

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

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

$$p(\mathbf{x}|y, \boldsymbol{\theta}) = \prod_{j=1}^D p(x_j|y=c, \boldsymbol{\theta}_{jc}) \quad \text{where } \mathbf{x} \text{ is a } D\text{-dimensional feature vector} \quad (3.34)$$

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

$$p(\mathbf{x}_i, y_i|\boldsymbol{\theta}) = p(y_i|\boldsymbol{\pi}) \prod_j p(x_{ij}|\boldsymbol{\theta}_j) = \prod_c \pi_c^{\mathbb{I}(y_i=c)} \prod_j \prod_c p(x_{ij}|\boldsymbol{\theta}_{jc})^{\mathbb{I}(y_i=c)} \quad (3.35)$$

Hence the log-likelihood is given by

$$p(\mathcal{D}|\boldsymbol{\theta}) = \sum_{c=1}^C N_c \log \pi_c + \sum_{j=1}^D \sum_{c=1}^C \sum_{i:y_i=c} \log p(x_{ij}|\boldsymbol{\theta}_{jc}) \quad (3.36)$$

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 (3.30), the MLE for the class prior is given by

$$\hat{\pi}_c = \frac{N_c}{N} \quad (3.37)$$

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

$$\hat{\boldsymbol{\theta}}_{jc} = \frac{N_{jc}}{N_c} \quad (3.38)$$

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

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

$$\hat{\theta}_{jc} = \left( \frac{N_{j1c}}{N_c}, \frac{N_{j2c}}{N_c}, \dots, \frac{N_{js_jc}}{N_c} \right)^T \quad (3.39)$$

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

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

## 3.5.2 Using the model for prediction

The goal is to compute

$$y = f(\mathbf{x}) = \arg \max_c P(y = c | \mathbf{x}, \boldsymbol{\theta}) = P(y = c | \boldsymbol{\theta}) \prod_{j=1}^D P(x_j | y = c, \boldsymbol{\theta}) \quad (3.40)$$

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

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:

$$\log p(y = c | \mathbf{x}, \boldsymbol{\theta}) = b_c - \log \left( \sum_{c'} e^{b_{c'}} \right) \quad \text{where } b_c \triangleq \log p(\mathbf{x} | y = c, \boldsymbol{\theta}) + \log p(y = c | \boldsymbol{\theta}) \quad (3.41)$$

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

$$\log(e^{-120} + e^{-121}) = \log(e^{-120}(1 + e^{-1})) = \log(1 + e^{-1}) - 120 \quad (3.42)$$

In general, we have

$$\sum_c e^{b_c} = \log \left[ \left( \sum_c e^{b_c - B} \right) e^B \right] = \log \left( \sum_c e^{b_c - B} \right) + B \quad \text{where } B \triangleq \max \{b_c\} \quad (3.43)$$

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 2.8.3) between feature  $X_j$  and the class label  $Y$

$$\mathbb{I}(X_j, Y) = \sum_{x_j} \sum_y p(x_j, y) \log \frac{p(x_j, y)}{p(x_j)p(y)} \quad (3.44)$$

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

$$\mathbb{I}_j = \sum_c \left[ \theta_{jc} \pi_c \log \frac{\theta_{jc}}{\theta_j} + (1 - \theta_{jc}) \pi_c \log \frac{1 - \theta_{jc}}{1 - \theta_j} \right] \quad (3.45)$$

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. 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:

$$p(\mathbf{x}_i | y_i = c, \boldsymbol{\theta}) = \quad (3.46)$$

#### 3.5.5.1 Bernoulli product model

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:

$$p(\mathbf{x}_i | y_i = c, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Ber}(x_{ij} | \theta_{jc}) = \prod_{j=1}^D \theta_{jc}^{\mathbb{I}(x_{ij})} (1 - \theta_{jc})^{1 - \mathbb{I}(x_{ij})} \quad (3.47)$$

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:

$$p(\mathbf{x}_i | y_i = c, \boldsymbol{\theta}) = \text{Mu}(\mathbf{x}_i | N_i, \boldsymbol{\theta}_c) = \frac{N_i!}{\prod_{j=1}^D x_{ij}!} \prod_{j=1}^D \theta_{jc}^{x_{ij}} \quad (3.48)$$

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 (3.48) 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:

$$p(\mathbf{x}_i | y_i = c, \boldsymbol{\alpha}) = \int \text{Mu}(\mathbf{x}_i | N_i, \boldsymbol{\theta}_c) \text{Dir}(\boldsymbol{\theta}_c | \boldsymbol{\alpha}_c) = \frac{N_i!}{\prod_{j=1}^D x_{ij}!} \prod_{j=1}^D \frac{B(\mathbf{x}_i + \boldsymbol{\alpha}_c)}{B(\boldsymbol{\alpha}_c)} \quad (3.49)$$

(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  $j$ , the posterior counts on  $j$  gets updated, making another occurrence of word  $j$  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

In this chapter, we 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 we will encounter in later chapters.

#### 4.1 Basics

Recall from Section 2.5.2 that the pdf for an MVN in  $D$  dimensions is defined by the following:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (4.1)$$

The expression inside the exponent is the Mahalanobis distance between a data vector  $\mathbf{x}$  and the mean vector  $\boldsymbol{\mu}$ . We can gain a better understanding of this quantity by performing an **eigendecomposition** of  $\boldsymbol{\Sigma}$ . That is, we write  $\boldsymbol{\Sigma} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T$ , where  $\mathbf{U}$  is an orthonormal matrix of eigenvectors satisfying  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ , and  $\boldsymbol{\Lambda}$  is a diagonal matrix of eigenvalues. Using the eigendecomposition, we have that

$$\boldsymbol{\Sigma}^{-1} = \mathbf{U}^{-T} \boldsymbol{\Lambda}^{-1} \mathbf{U}^{-1} = \mathbf{U} \boldsymbol{\Lambda}^{-1} \mathbf{U}^T = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad (4.2)$$

where  $\mathbf{u}_i$  is the  $i$ th column of  $\mathbf{U}$ , containing the  $i$ th eigenvector. Hence we can rewrite the Mahalanobis distance as follows:

$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^T \left( \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \right) (\mathbf{x} - \boldsymbol{\mu}) \quad (4.3)$$

$$= \sum_{i=1}^D \frac{1}{\lambda_i} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{u}_i \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \quad (4.4)$$

where  $y_i \triangleq \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})$ . Recall that the equation for an ellipse in 2d is

$$\frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2} = 1 \quad (4.5)$$

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.

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

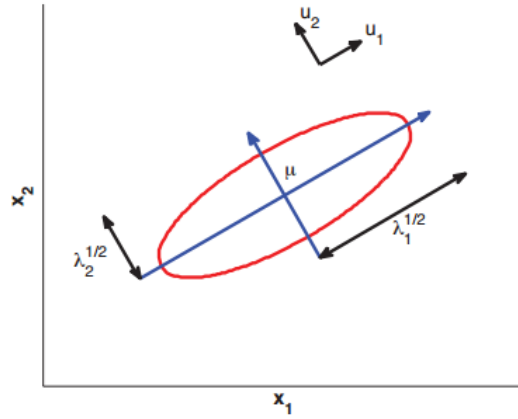


Fig. 4.1: Visualization of a 2 dimensional Gaussian density. The major and minor axes of the ellipse are defined by the first two eigenvectors of the covariance matrix, namely  $u_1$  and  $u_2$ . Based on Figure 2.7 of (Bishop 2006a)

#### 4.1.1 MLE for a MVN

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

$$\bar{\mu} = \frac{1}{N} \sum_{i=1}^N x_i \triangleq \bar{x} \quad (4.6)$$

$$\bar{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T = \frac{1}{N} \left( \sum_{i=1}^N x_i x_i^T \right) - \bar{x} \bar{x}^T \quad (4.7)$$

#### 4.1.2 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

$$f(x) = \frac{1}{Z} \exp \left( -\frac{1}{2} x^T \Sigma^{-1} x \right) \quad (4.8)$$

### 4.2 Gaussian discriminant analysis

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

$$p(x|y=c, \theta) = \mathcal{N}(x|\mu_c, \Sigma_c) \quad (4.9)$$

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  $\Sigma_c$  is diagonal, this is equivalent to naive Bayes.

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

$$y = \arg \max_c [\log p(y = c | \pi) + \log p(\mathbf{x} | \theta)] \quad (4.10)$$

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,  $\mu_c$ , using Mahalanobis distance. This can be thought of as a **nearest centroids classifier**.

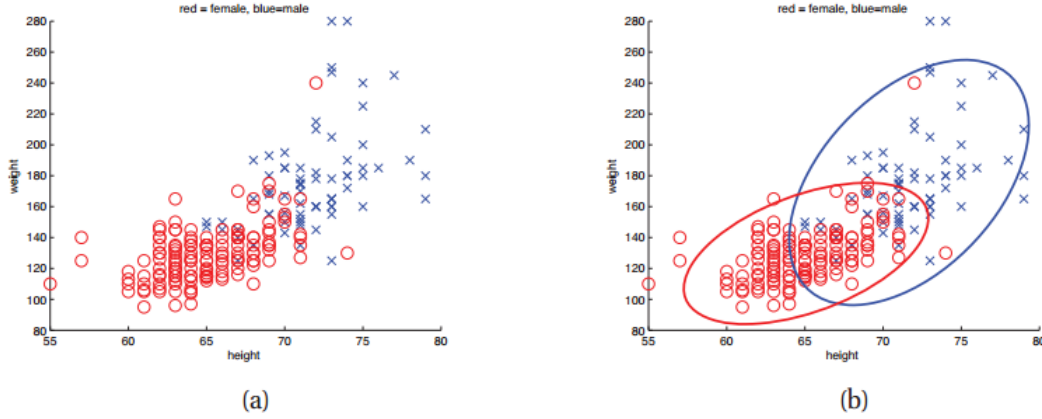


Fig. 4.2: (a) Height/weight data. (b) Visualization of 2d Gaussians fit to each class. 95% of the probability mass is inside the ellipse.

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:

$$y = \arg \max_c (\mathbf{x} - \mu_c)^T \Sigma_c^{-1} (\mathbf{x} - \mu_c) \quad (4.11)$$

### 4.2.1 Quadratic discriminant analysis (QDA)

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

$$p(y | \mathbf{x}, \theta) = \frac{\pi_c |2\pi \Sigma_c|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma_c^{-1} (\mathbf{x} - \mu) \right]}{\sum_{c'} \pi_{c'} |2\pi \Sigma_{c'}|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma_{c'}^{-1} (\mathbf{x} - \mu) \right]} \quad (4.12)$$

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,  $\Sigma_c = \Sigma$ . In this case, we can simplify Equation (4.12) as follows:

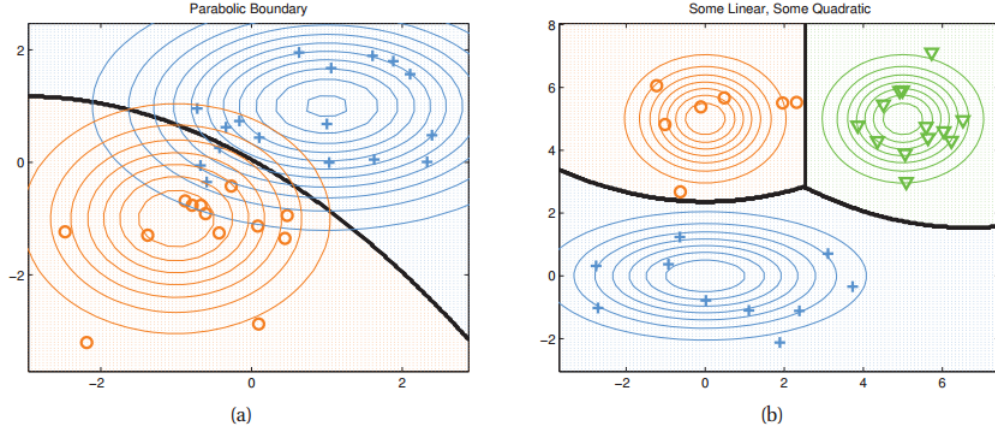


Fig. 4.3: Quadratic decision boundaries in 2D for the 2 and 3 class case.

$$p(y|\mathbf{x}, \boldsymbol{\theta}) \propto \pi_c \exp \left( \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c \right) \quad (4.13)$$

$$= \exp \left( \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \right) \exp \left( -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} \right) \quad (4.14)$$

$$\propto \exp \left( \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \right) \quad (4.15)$$

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

$$\gamma_c \triangleq -\frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \quad (4.16)$$

$$\boldsymbol{\beta}_c \triangleq \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c \quad (4.17)$$

then we can write

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\boldsymbol{\beta}_c^T \mathbf{x} + \gamma_c}}{\sum_{c'} e^{\boldsymbol{\beta}_{c'}^T \mathbf{x} + \gamma_{c'}}} \triangleq \sigma(\boldsymbol{\eta}, c) \quad \text{where } \boldsymbol{\eta} \triangleq (e^{\boldsymbol{\beta}_1^T \mathbf{x} + \gamma_1}, \dots, e^{\boldsymbol{\beta}_C^T \mathbf{x} + \gamma_C}) \quad (4.18)$$

where  $\sigma()$  is the **softmax activation function**<sup>16</sup>, defined as follows:

$$\sigma(\mathbf{q}, i) \triangleq \frac{\exp(q_i)}{\sum_{j=1}^n \exp(q_j)} \quad (4.19)$$

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

$$\mathcal{S}_\alpha(\mathbf{x}) = \frac{\sum_{j=1}^D x_j e^{\alpha x_j}}{\sum_{j=1}^D e^{\alpha x_j}} \quad (4.20)$$

$\mathcal{S}_\alpha$  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.

<sup>16</sup> [http://en.wikipedia.org/wiki/Softmax\\_activation\\_function](http://en.wikipedia.org/wiki/Softmax_activation_function)

An interesting property of Equation (4.18) is that, if we take logs, we end up with a linear function of  $\mathbf{x}$ . (The reason it is linear is because the  $\mathbf{x}^T \Sigma^{-1} \mathbf{x}$  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

$$p(y = 1|\mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\beta_1^T \mathbf{x} + \gamma_1}}{e^{\beta_0^T \mathbf{x} + \gamma_0} + e^{\beta_1^T \mathbf{x} + \gamma_1}} \quad (4.21)$$

$$= \frac{1}{1 + e^{(\beta_0 - \beta_1)^T \mathbf{x} + (\gamma_0 - \gamma_1)}} = \text{sigm}((\beta_1 - \beta_0)^T \mathbf{x} + (\gamma_0 - \gamma_1)) \quad (4.22)$$

where  $\text{sigm}(x)$  refers to the sigmoid function<sup>17</sup>.

Now

$$\gamma_1 - \gamma_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_0^T \Sigma^{-1} \boldsymbol{\mu}_0 + \log(\pi_1/\pi_0) \quad (4.23)$$

$$= -\frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^T \Sigma^{-1} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) + \log(\pi_1/\pi_0) \quad (4.24)$$

So if we define

$$\mathbf{w} = \beta_1 - \beta_0 = \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \quad (4.25)$$

$$\mathbf{x}_0 = \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) - (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \frac{\log(\pi_1/\pi_0)}{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^T \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)} \quad (4.26)$$

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

$$p(y = 1|\mathbf{x}, \boldsymbol{\theta}) = \text{sigm}(\mathbf{w}^T (\mathbf{x} - \mathbf{x}_0)) \quad (4.27)$$

(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.

If  $\Sigma = \sigma^2 \mathbf{I}$ , then  $\mathbf{w}$  is in the direction of  $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0$ . So we classify the point based on whether its projection is closer to  $\boldsymbol{\mu}_0$  or  $\boldsymbol{\mu}_1$ . This is illustrated in Figure 4.4. Furthermore, if  $\pi_1 = \pi_0$ , then  $\mathbf{x}_0 = \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0)$ , which is half way between the means. If we make  $\pi_1 > \pi_0$ , then  $\mathbf{x}_0$  gets closer to  $\boldsymbol{\mu}_0$ , 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**:

$$d' \triangleq \frac{\mu_1 - \mu_0}{\sigma} \quad (4.28)$$

<sup>17</sup> [http://en.wikipedia.org/wiki/Sigmoid\\_function](http://en.wikipedia.org/wiki/Sigmoid_function)

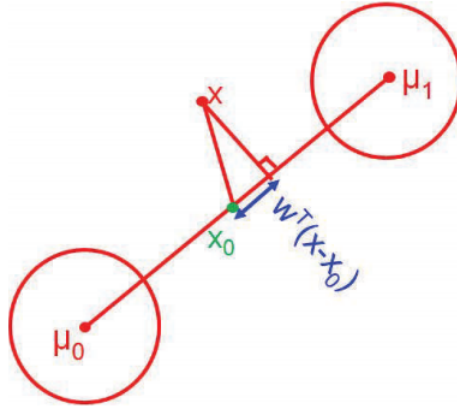


Fig. 4.4: Geometry of LDA in the 2 class case where  $\Sigma_1 = \Sigma_2 = I$ .

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:

$$p(\mathcal{D}|\theta) = \sum_{c=1}^C \sum_{i:y_i=c} \log \pi_c + \sum_{c=1}^C \sum_{i:y_i=c} \log \mathcal{N}(x_i|\mu_c, \Sigma_c) \quad (4.29)$$

The MLE for each parameter is as follows:

$$\bar{\mu}_c = \frac{N_c}{N} \quad (4.30)$$

$$\bar{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} x_i \quad (4.31)$$

$$\bar{\Sigma}_c = \frac{1}{N_c} \sum_{i:y_i=c} (x_i - \bar{\mu}_c)(x_i - \bar{\mu}_c)^T \quad (4.32)$$

#### 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 kinds of 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

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}, \boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1} = \begin{pmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{pmatrix}, \quad (4.33)$$

Then the marginals are given by

$$\begin{aligned} p(\mathbf{x}_1) &= \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}) \\ p(\mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_2 | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}) \end{aligned} \quad (4.34)$$

and the posterior conditional is given by

$$\begin{aligned} p(\mathbf{x}_1|\mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2}) \\ \boldsymbol{\mu}_{1|2} &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{11}^{-1} \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\Sigma}_{1|2} (\boldsymbol{\Lambda}_{11} \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2)) \\ \boldsymbol{\Sigma}_{1|2} &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} = \boldsymbol{\Lambda}_{11}^{-1} \end{aligned} \quad (4.35)$$

Equation (4.35) 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  $x_1$  or  $x_2$ . 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  $x_2$ , and the conditional covariance is just a constant matrix that is independent of  $x_2$ . 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

## 4.4 Inferring the parameters of an MVN

### 4.4.1 Posterior distribution of $\mu$

### 4.4.2 Posterior distribution of $\Sigma$ \*

### 4.4.3 Posterior distribution of $\mu$ and $\Sigma$ \*

### 4.4.4 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.

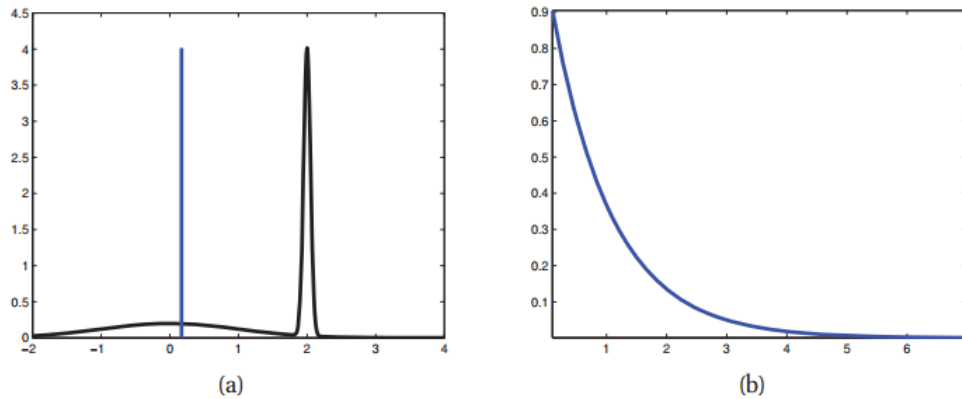


Fig. 5.1: (a) A bimodal distribution in which the mode is very untypical of the distribution. The thin blue vertical line is the mean, which is arguably a better summary of the distribution, since it is near the majority of the probability mass. (b) A skewed distribution in which the mode is quite different 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 1.2.2.1), 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 (2.43). 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 2.7). 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.

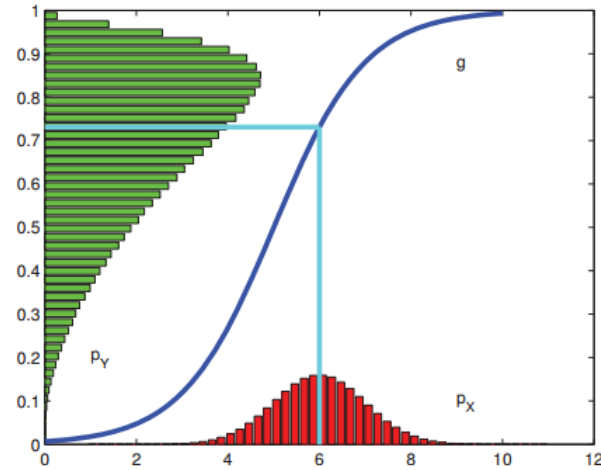


Fig. 5.2: Example of the transformation of a density under a nonlinear transform. Note how the mode of the transformed distribution is not the transform of the original mode. Based on Exercise 1.4 of (Bishop 2006b).

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.

### 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>18</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.

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 (2.3)):

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

<sup>18</sup> This example is from <http://www.johndcook.com/blog/2011/09/27/bayesian-amazon/>

This is similar to a leave-one-out cross-validation estimate (Section 1.3) 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.

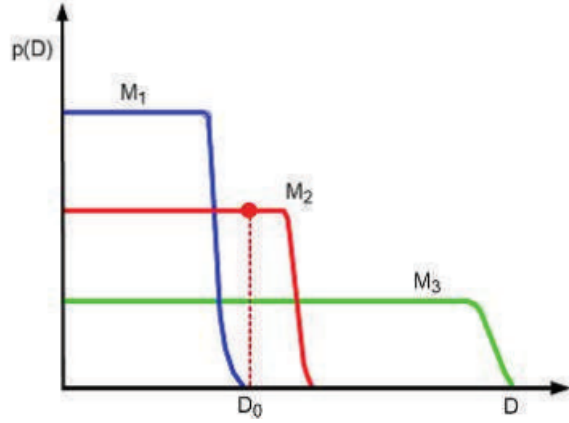


Fig. 5.3: A schematic illustration of the Bayesian Occams razor. The broad (green) curve corresponds to a complex model, the narrow (blue) curve to a simple model, and the middle (red) curve is just right. Based on Figure 3.13 of (Bishop 2006a).

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



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

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

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

#### 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 (1.1)).

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**:

$$\rho(f) = \mathbb{E}_{p(y|x)}[L(y, f(x))] = \begin{cases} \sum_y L[y, f(x)]p(y|x) & , y \text{ is discrete} \\ \int_y L[y, f(x)]p(y|x)dy & , y \text{ is continuous} \end{cases} \quad (5.28)$$

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 1.2.2.1), we can prove 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

$$\rho(f) = \int_y L[y, f(\mathbf{x})] p(y | \mathbf{x}) dy = \int_y [y - f(\mathbf{x})]^2 p(y | \mathbf{x}) dy = \int_y [y^2 - 2yf(\mathbf{x}) + f(\mathbf{x})^2] p(y | \mathbf{x}) dy \quad (5.30)$$

Hence the optimal estimate is the posterior mean:

$$\begin{aligned}
\frac{\partial \rho}{\partial f} &= \int_y [-2y + 2f(\mathbf{x})] p(y | \mathbf{x}) dy = 0 \Rightarrow \\
\int_y f(\mathbf{x}) p(y | \mathbf{x}) dy &= \int_y y p(y | \mathbf{x}) dy \\
f(\mathbf{x}) \int_y p(y | \mathbf{x}) dy &= \mathbb{E}_{p(y | \mathbf{x})} [y] \\
f(\mathbf{x}) &= \mathbb{E}_{p(y | \mathbf{x})} [y] \quad (5.31)
\end{aligned}$$

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.*

$$\begin{aligned}
\rho(f) &= \int_y L[y, f(\mathbf{x})] p(y|\mathbf{x}) dy = \int_y |y - f(\mathbf{x})| p(y|\mathbf{x}) dy \\
&= \int_y \{ [f(\mathbf{x}) - y] p(y < f(\mathbf{x})|\mathbf{x}) + [y - f(\mathbf{x})] p(y \geq f(\mathbf{x})|\mathbf{x}) \} dy \\
\frac{\partial \rho}{\partial f} &= \int_y [p(y < f(\mathbf{x})|\mathbf{x}) - p(y \geq f(\mathbf{x})|\mathbf{x})] dy = 0 \Rightarrow \\
p(y < f(\mathbf{x})|\mathbf{x}) &= p(y \geq f(\mathbf{x})|\mathbf{x}) = 0.5 \\
\therefore f(\mathbf{x}) &= \text{median}
\end{aligned}$$

#### 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 don't 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

### Frequentist statistics

Attempts have been made to devise approaches to statistical inference that avoid treating parameters like random variables, and which thus avoid the use of priors and Bayes rule. Such approaches are known as **frequentist statistics**, **classical statistics** or **orthodox statistics**. Instead of being based on the posterior distribution, they are based on the concept of a sampling distribution.

#### 6.1 Sampling distribution of an estimator

In frequentist statistics, a parameter estimate  $\hat{\theta}$  is computed by applying an **estimator**  $\delta$  to some data  $\mathcal{D}$ , so  $\hat{\theta} = \delta(\mathcal{D})$ . The parameter is viewed as fixed and the data as random, which is the exact opposite of the Bayesian approach. The uncertainty in the parameter estimate can be measured by computing the **sampling distribution** of the estimator. To understand this

##### 6.1.1 Bootstrap

We might think of the bootstrap distribution as a poor mans Bayes posterior, see (Hastie et al. 2001, p235) for details.

##### 6.1.2 Large sample theory for the MLE \*

#### 6.2 Frequentist decision theory

In frequentist or classical decision theory, there is a loss function and a likelihood, but there is no prior and hence no posterior or posterior expected loss. Thus there is no automatic way of deriving an optimal estimator, unlike the Bayesian case. Instead, in the frequentist approach, we are free to choose any estimator or decision procedure  $f : \mathcal{X} \rightarrow \mathcal{Y}$  we want.

Having chosen an estimator, we define its expected loss or **risk** as follows:

$$R_{\text{exp}}(\theta, f) \triangleq \mathbb{E}_{p(\tilde{\mathcal{D}}|\theta^*)}[L(\theta^*, f(\tilde{\mathcal{D}}))] = \int L(\theta^*, f(\tilde{\mathcal{D}}))p(\tilde{\mathcal{D}}|\theta^*)d\tilde{\mathcal{D}} \quad (6.1)$$

where  $\tilde{\mathcal{D}}$  is data sampled from nature's distribution, which is represented by parameter  $\theta^*$ . In other words, the expectation is wrt the sampling distribution of the estimator. Compare this to the Bayesian posterior expected loss:

$$\rho(f|\mathcal{D},) \quad (6.2)$$

### **6.3 Desirable properties of estimators**

### **6.4 Empirical risk minimization**

#### **6.4.1 Regularized risk minimization**

#### **6.4.2 Structural risk minimization**

#### **6.4.3 Estimating the risk using cross validation**

#### **6.4.4 Upper bounding the risk using statistical learning theory \***

#### **6.4.5 Surrogate loss functions**

**log-loss**

$$L_{\text{nl}}(y, \eta) = -\log p(y|x, w) = \log(1 + e^{-y\eta}) \quad (6.3)$$

### **6.5 Pathologies of frequentist statistics \***



## Chapter 7

# Linear Regression

### 7.1 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.

### 7.2 Representation

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mathbf{w}^T \mathbf{x}, \sigma^2) \quad (7.1)$$

where  $\mathbf{w}$  and  $\mathbf{x}$  are extended vectors,  $\mathbf{x} = (1, x)$ ,  $\mathbf{w} = (b, w)$ .

Linear regression can be made to model non-linear relationships by replacing  $\mathbf{x}$  with some non-linear function of the inputs,  $\phi(\mathbf{x})$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mathbf{w}^T \phi(\mathbf{x}), \sigma^2) \quad (7.2)$$

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

$$\phi(x) = (1, x, \dots, x^d) \quad (7.3)$$

### 7.3 MLE

Instead of maximizing the log-likelihood, we can equivalently minimize the **negative log likelihood** or **NLL**:

$$\text{NLL}(\boldsymbol{\theta}) \triangleq -\ell(\boldsymbol{\theta}) = -\log(\mathcal{D}|\boldsymbol{\theta}) \quad (7.4)$$

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

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^N \log \left[ \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{1}{2\sigma^2} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \right) \right] \quad (7.5)$$

$$= -\frac{1}{2\sigma^2} \text{RSS}(\mathbf{w}) - \frac{N}{2} \log(2\pi\sigma^2) \quad (7.6)$$

RSS stands for **residual sum of squares** and is defined by

$$\text{RSS}(\mathbf{w}) \triangleq \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \quad (7.7)$$

We see that the MLE for  $\mathbf{w}$  is the one that minimizes the RSS, so this method is known as **least squares**. Let's drop constants wrt  $\mathbf{w}$  and NLL can be written as

$$\text{NLL}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \quad (7.8)$$

There two ways to minimize  $\text{NLL}(\mathbf{w})$ .

### 7.3.1 OLS

Define  $\mathbf{y} = (y_1, y_2, \dots, y_N)$ ,  $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix}$ , then  $\text{NLL}(\mathbf{w})$  can be written as

$$\text{NLL}(\mathbf{w}) = \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \quad (7.9)$$

When  $\mathcal{D}$  is small (for example,  $N < 1000$ ), we can use the following equation to compute  $\mathbf{w}$  directly

$$\hat{\mathbf{w}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (7.10)$$

The corresponding solution  $\hat{\mathbf{w}}_{\text{OLS}}$  to this linear system of equations is called the **ordinary least squares** or **OLS** solution.

*Proof.* We now state without proof some facts of matrix derivatives (we wont need all of these at this section).

$$\text{tr} A \triangleq \sum_{i=1}^n A_{ii}$$

$$\frac{\partial}{\partial A} AB = B^T \quad (7.11)$$

$$\frac{\partial}{\partial A^T} f(A) = \left[ \frac{\partial}{\partial A} f(A) \right]^T \quad (7.12)$$

$$\frac{\partial}{\partial A} ABA^T C = CAB + C^T AB^T \quad (7.13)$$

$$\frac{\partial}{\partial A} |A| = |A| (A^{-1})^T \quad (7.14)$$

Then,

$$\begin{aligned}
\text{NLL}(\mathbf{w}) &= \frac{1}{2N}(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) \\
\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y}) \\
&= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w}) \\
&= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} \text{tr}(\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w}) \\
&= \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} (\text{tr} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2 \text{tr} \mathbf{y}^T \mathbf{X} \mathbf{w})
\end{aligned}$$

Combining Equations (7.12) and (7.13), 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,

$$\begin{aligned}
\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= \frac{1}{2} (\mathbf{X}^T \mathbf{X} \mathbf{w} + \mathbf{X}^T \mathbf{X} \mathbf{w} - 2 \mathbf{X}^T \mathbf{y}) \\
&= \frac{1}{2} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) \\
\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= 0 \Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = 0 \\
\mathbf{X}^T \mathbf{X} \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\
\hat{\mathbf{w}}_{\text{OLS}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}
\end{aligned} \tag{7.15}$$

Equation (7.15) is known as the **normal equation**.

### 7.3.1.1 Geometric interpretation

See Figure 7.1.

To minimize the norm of the residual,  $\mathbf{y} - \hat{\mathbf{y}}$ , we want the residual vector to be orthogonal to every column of  $\mathbf{X}$ , so  $\tilde{\mathbf{x}}_j(\mathbf{y} - \hat{\mathbf{y}}) = 0$  for  $j = 1 : D$ . Hence

$$\tilde{\mathbf{x}}_j(\mathbf{y} - \hat{\mathbf{y}}) = 0 \Rightarrow \mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) = 0 \Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \tag{7.16}$$

### 7.3.2 SGD

When  $\mathcal{D}$  is large, use stochastic gradient descent(SGD).

$$\therefore \frac{\partial}{\partial w_i} \text{NLL}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i) x_{ij} \tag{7.17}$$

$$\begin{aligned}
\therefore w_j &= w_j - \alpha \frac{\partial}{\partial w_j} \text{NLL}(\mathbf{w}) \\
&= w_j - \sum_{i=1}^N \alpha (\mathbf{w}^T \mathbf{x}_i - y_i) x_{ij}
\end{aligned} \tag{7.18}$$

$$\therefore \mathbf{w} = \mathbf{w} - \alpha (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x} \tag{7.19}$$

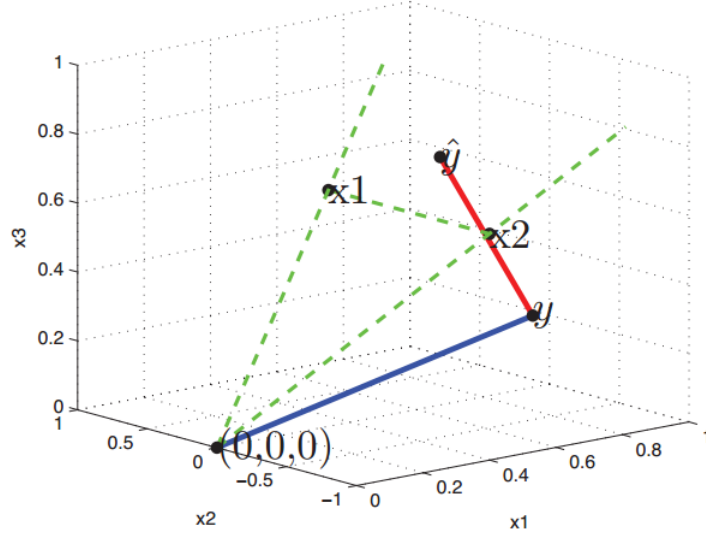


Fig. 7.1: Graphical interpretation of least squares for  $N = 3$  examples and  $D = 2$  features.  $\tilde{x}_1$  and  $\tilde{x}_2$  are vectors in  $\mathbb{R}^3$ ; together they define a 2D plane.  $y$  is also a vector in  $\mathbb{R}^3$  but does not lie on this 2D plane. The orthogonal projection 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.

## 7.4 Ridge regression(MAP)

One problem with ML estimation is that it can result in overfitting. In this section, we discuss a way to ameliorate this problem by using MAP estimation with a Gaussian prior.

### 7.4.1 Basic idea

We can encourage the parameters to be small, thus resulting in a smoother curve, by using a zero-mean Gaussian prior:

$$p(w) = \prod_j \mathcal{N}(w_j | 0, \tau^2) \quad (7.20)$$

where  $1/\tau^2$  controls the strength of the prior. The corresponding MAP estimation problem becomes

$$\arg \max_w \sum_{i=1}^N \log \mathcal{N}(y_i | w_0 + w^T x_i, \sigma^2) + \sum_{j=1}^D \log \mathcal{N}(w_j | 0, \tau^2) \quad (7.21)$$

It is a simple exercise to show that this is equivalent to minimizing the following

$$J(w) = \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + w^T x_i))^2 + \lambda \|w\|^2, \text{ where } \lambda \triangleq \frac{\sigma^2}{\tau^2} \quad (7.22)$$

Here the first term is the MSE/NLL as usual, and the second term,  $\lambda \geq 0$ , is a complexity penalty. The corresponding solution is given by

$$\hat{w}_{\text{ridge}} = (\lambda I_D + X^T X)^{-1} X^T y \quad (7.23)$$

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.

### 7.4.2 Numerically stable computation \*

$$\hat{\mathbf{w}}_{\text{ridge}} = \mathbf{V}(\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I}_N)^{-1} \mathbf{Z}^T \mathbf{y} \quad (7.24)$$

### 7.4.3 Connection with PCA \*

### 7.4.4 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).

## 7.5 Bayesian linear regression

TODO



## Chapter 8

# Logistic Regression

### 8.1 Representation

Logistic regression can be binomial or multinomial. The **binomial logistic regression** model has the following form

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Ber}(y|\text{sigm}(\mathbf{w}^T \mathbf{x})) \quad (8.1)$$

where  $\mathbf{w}$  and  $\mathbf{x}$  are extended vectors, i.e.,  $\mathbf{w} = (b, w_1, w_2, \dots, w_D)$ ,  $\mathbf{x} = (1, x_1, x_2, \dots, x_D)$ .

### 8.2 Optimization

#### 8.2.1 MLE

$$\begin{aligned} \ell(\mathbf{w}) &= \log \left\{ \prod_{i=1}^N [\pi(\mathbf{x}_i)]^{y_i} [1 - \pi(\mathbf{x}_i)]^{1-y_i} \right\}, \text{ where } \pi(\mathbf{x}) \triangleq P(y = 1|\mathbf{x}, \mathbf{w}) \\ &= \sum_{i=1}^N [y_i \log \pi(\mathbf{x}_i) + (1 - y_i) \log(1 - \pi(\mathbf{x}_i))] \end{aligned} \quad (8.2)$$

$$\begin{aligned} &= \sum_{i=1}^N \left[ y_i \log \frac{\pi(\mathbf{x}_i)}{1 - \pi(\mathbf{x}_i)} + \log(1 - \pi(\mathbf{x}_i)) \right] \\ &= \sum_{i=1}^N [y_i(\mathbf{w} \cdot \mathbf{x}_i) - \log(1 + \exp(\mathbf{w} \cdot \mathbf{x}_i))] \\ J(\mathbf{w}) &\triangleq \text{NLL}(\mathbf{w}) = -\ell(\mathbf{w}) = -\sum_{i=1}^N [y_i(\mathbf{w} \cdot \mathbf{x}_i) - \log(1 + \exp(\mathbf{w} \cdot \mathbf{x}_i))] \end{aligned} \quad (8.3)$$

Equation (8.2) is also called the **cross-entropy** error function (see Equation (2.50)).

Unlike linear regression, we can no longer write down the MLE in closed form. Instead, we need to use an optimization algorithm to compute it, see Appendix A. For this, we need to derive the gradient and Hessian.

In the case of logistic regression, one can show that the gradient and Hessian of this are given by the following

$$\mathbf{g}(\mathbf{w}) = \frac{dJ}{d\mathbf{w}} = \sum_{i=1}^N [\pi(\mathbf{x}_i) - y_i] \mathbf{x}_i = \mathbf{X}(\boldsymbol{\pi} - \mathbf{y}) \quad (8.4)$$

$$\begin{aligned} \mathbf{H}(\mathbf{w}) &= \frac{d\mathbf{g}^T}{d\mathbf{w}} = \frac{d}{d\mathbf{w}} (\boldsymbol{\pi} - \mathbf{y})^T \mathbf{X}^T \\ &= \frac{d}{d\mathbf{w}} \boldsymbol{\pi}^T \mathbf{X}^T \\ &= (\pi(\mathbf{x}_1)(1 - \pi(\mathbf{x}_1))\mathbf{x}_1, \dots) \mathbf{X}^T \\ &= \mathbf{X} \mathbf{S} \mathbf{X}^T, \text{ where } \mathbf{S} \triangleq \text{diag}(\pi(\mathbf{x}_1)(1 - \pi(\mathbf{x}_1))\mathbf{x}_1) \end{aligned} \quad (8.5)$$

### 8.2.1.1 Iteratively reweighted least squares (IRLS)

TODO

## 8.2.2 MAP

Just as we prefer ridge regression to linear regression, so we should prefer MAP estimation for logistic regression to computing the MLE.  $\ell_2$  regularization

we can use  $\ell_2$  regularization, just as we did with ridge regression. We note that the new objective, gradient and Hessian have the following forms:

$$J'(\mathbf{w}) \triangleq \text{NLL}(\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w} \quad (8.6)$$

$$\mathbf{g}'(\mathbf{w}) = \mathbf{g}(\mathbf{w}) + \lambda \mathbf{w} \quad (8.7)$$

$$\mathbf{H}'(\mathbf{w}) = \mathbf{H}(\mathbf{w}) + \lambda \mathbf{I} \quad (8.8)$$

$$(8.9)$$

It is a simple matter to pass these modified equations into any gradient-based optimizer.

## 8.3 Multinomial logistic regression

### 8.3.1 Representation

**Multinomial logistic regression** model is also called a **maximum entropy classifier**, which has the following form

$$p(y = c | \mathbf{x}, \mathbf{W}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{c=1}^C \exp(\mathbf{w}_c^T \mathbf{x})} \quad (8.10)$$

### 8.3.2 MLE

Let  $\mathbf{y}_i = (\mathbb{I}(y_i = 1), \mathbb{I}(y_i = 2), \dots, \mathbb{I}(y_i = C))$ ,  $\boldsymbol{\mu}_i = (p(y = 1 | \mathbf{x}_i, \mathbf{W}), p(y = 2 | \mathbf{x}_i, \mathbf{W}), \dots, p(y = C | \mathbf{x}_i, \mathbf{W}))$ , then the log-likelihood function can be written as



$$\ell(\mathbf{W}) = \log \prod_{i=1}^N \prod_{c=1}^C \mu_{ic}^{y_{ic}} = \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log \mu_{ic} \quad (8.11)$$

$$= \sum_{i=1}^N \left[ \left( \sum_{c=1}^C y_{ic} \mathbf{w}_c^T \mathbf{x}_i \right) - \log \left( \sum_{c=1}^C \exp(\mathbf{w}_c^T \mathbf{x}_i) \right) \right] \quad (8.12)$$

Define the objective function as NLL

$$J(\mathbf{W}) = \text{NLL}(\mathbf{W}) = -\ell(\mathbf{W}) \quad (8.13)$$

Define  $\mathbf{A} \otimes \mathbf{B}$  be the **kronecker product** of matrices  $\mathbf{A}$  and  $\mathbf{B}$ . If  $\mathbf{A}$  is an  $m \times n$  matrix and  $\mathbf{B}$  is a  $p \times q$  matrix, then  $\mathbf{A} \otimes \mathbf{B}$  is the  $mp \times nq$  block matrix

$$\mathbf{A} \otimes \mathbf{B} \triangleq \begin{pmatrix} a_{11}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{pmatrix} \quad (8.14)$$

The gradient and Hessian are given by

$$\mathbf{g}(\mathbf{W}) = \sum_{i=1}^N (\boldsymbol{\mu} - \mathbf{y}_i) \otimes \mathbf{x}_i \quad (8.15)$$

$$\mathbf{H}(\mathbf{W}) = \sum_{i=1}^N (\text{diag}(\boldsymbol{\mu}_i) - \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T) \otimes (\mathbf{x}_i \mathbf{x}_i^T) \quad (8.16)$$

$$(8.17)$$

where  $\mathbf{y}_i = (\mathbb{I}(y_i = 1), \mathbb{I}(y_i = 1), \dots, \mathbb{I}(y_i = C - 1))$  and  $\boldsymbol{\mu}_i = (p(y = 1 | \mathbf{x}_i, \mathbf{W}), p(y = 2 | \mathbf{x}_i, \mathbf{W}), \dots, p(y = C - 1 | \mathbf{x}_i, \mathbf{W}))$  are column vectors of length  $C - 1$ .

Pass them to any gradient-based optimizer.

### 8.3.3 MAP

The new objective

$$J'(\mathbf{W}) = \text{NLL}(\mathbf{w}) - \log p(\mathbf{W}) \quad , \text{ where } p(\mathbf{W}) \triangleq \prod_{c=1}^C \mathcal{N}(\mathbf{w}_c | \mathbf{0}, \mathbf{V}_0) \quad (8.18)$$

$$= J(\mathbf{W}) + \frac{1}{2} \sum_{c=1}^C \mathbf{w}_c^T \mathbf{V}_0^{-1} \mathbf{w}_c \quad (8.19)$$

$$(8.20)$$

Its gradient and Hessian are given by

$$\mathbf{g}'(\mathbf{w}) = \mathbf{g}(\mathbf{W}) + \mathbf{V}_0^{-1} \left( \sum_{c=1}^C \mathbf{w}_c \right) \quad (8.21)$$

$$\mathbf{H}'(\mathbf{w}) = \mathbf{H}(\mathbf{w}) + \mathbf{I}_C \otimes \mathbf{V}_0^{-1} \quad (8.22)$$

This can be passed to any gradient-based optimizer to find the MAP estimate. Note, however, that the Hessian has size  $((CD)(CD))$ , which is  $C$  times more row and columns than in the binary case, so limited memory BFGS is more appropriate than Newtons method.

## 8.4 Bayesian logistic regression

It is natural to want to compute the full posterior over the parameters,  $p(\mathbf{w}|\mathcal{D})$ , for logistic regression models. This can be useful for any situation where we want to associate confidence intervals with our predictions (e.g., this is necessary when solving contextual bandit problems, discussed in Section TODO).

Unfortunately, unlike the linear regression case, this cannot be done exactly, since there is no convenient conjugate prior for logistic regression. We discuss one simple approximation below; some other approaches include MCMC (Section TODO), variational inference (Section TODO), expectation propagation (Kuss and Rasmussen 2005), etc. For notational simplicity, we stick to binary logistic regression.

### 8.4.1 Laplace approximation

### 8.4.2 Derivation of the BIC

### 8.4.3 Gaussian approximation for logistic regression

### 8.4.4 Approximating the posterior predictive

### 8.4.5 Residual analysis (outlier detection) \*

## 8.5 Online learning and stochastic optimization

Traditionally machine learning is performed **offline**, however, if we have **streaming data**, we need to perform **online learning**, so we can update our estimates as each new data point arrives rather than waiting until the end (which may never occur). And even if we have a batch of data, we might want to treat it like a stream if it is too large to hold in main memory. Below we discuss learning methods for this kind of scenario.

TODO

## 8.6 Generative vs discriminative classifiers

### 8.6.1 Pros and cons of each approach

- **Easy to fit?** As we have seen, it is usually very easy to fit generative classifiers. For example, in Sections 3.5.1 and 4.2.4, we show that we can fit a naive Bayes model and an LDA model by simple counting and averaging. By contrast, logistic regression requires solving a convex optimization problem (see Section 8.2 for the details), which is much slower.
- **Fit classes separately?** In a generative classifier, we estimate the parameters of each class conditional density independently, so we do not have to retrain the model when we add more classes. In contrast, in discriminative models, all the parameters interact, so the whole model must be retrained if we add a new class. (This is also the case if we train a generative model to maximize a discriminative objective Salojärvi et al. (2005).)
- **Handle missing features easily?** Sometimes some of the inputs (components of  $\mathbf{x}$ ) are not observed. In a generative classifier, there is a simple method for dealing with this, as we discuss in Section 8.6.2. However, in a discriminative classifier, there is no principled solution to this problem, since the model assumes that  $\mathbf{x}$  is always available to be conditioned on (although see (Marlin 2008) for some heuristic approaches).
- **Can handle unlabeled training data?** There is much interest in **semi-supervised learning**, which uses unlabeled data to help solve a supervised task. This is fairly easy to do using generative models (see e.g., (Lasserre et al. 2006; Liang et al. 2007)), but is much harder to do with discriminative models.

- **Symmetric in inputs and outputs?** We can run a generative model backwards, and infer probable inputs given the output by computing  $p(\mathbf{x}|\mathbf{y})$ . This is not possible with a discriminative model. The reason is that a generative model defines a joint distribution on  $\mathbf{x}$  and  $\mathbf{y}$ , and hence treats both inputs and outputs symmetrically.
- **Can handle feature preprocessing?** A big advantage of discriminative methods is that they allow us to preprocess the input in arbitrary ways, e.g., we can replace  $\mathbf{x}$  with  $\phi(\mathbf{x})$ , which could be some basis function expansion, as illustrated in Figure 8.1. It is often hard to define a generative model on such pre-processed data, since the new features are correlated in complex ways.
- **Well-calibrated probabilities?** Some generative models, such as naive Bayes, make strong independence assumptions which are often not valid. This can result in very extreme posterior class probabilities (very near 0 or 1). Discriminative models, such as logistic regression, are usually better calibrated in terms of their probability estimates.

See Table 8.1 for a summary of the classification and regression techniques we cover in this book.

Model	Classif/regr	Gen/Discr	Param/Non	Section
Discriminant analysis	Classif	Gen	Param	Sec. 4.2.2, 4.2.4
Naive Bayes classifier	Classif	Gen	Param	Sec. 3.5, 3.5.1.2
Tree-augmented Naive Bayes classifier	Classif	Gen	Param	Sec. 10.2.1
Linear regression	Regr	Discrim	Param	Sec. 1.4.5, 7.3, 7.6
Logistic regression	Classif	Discrim	Param	Sec. 1.4.6, 8.2.1.1, 8.4.3, 21.8.1.1
Sparse linear/ logistic regression	Both	Discrim	Param	Ch. 13
Mixture of experts	Both	Discrim	Param	Sec. 11.2.4
Multilayer perceptron (MLP)/ Neural network	Both	Discrim	Param	Ch. 16
Conditional random field (CRF)	Classif	Discrim	Param	Sec. 19.6
<i>K</i> nearest neighbor classifier	Classif	Gen	Non	Sec. TODO, TODO
(Infinite) Mixture Discriminant analysis	Classif	Gen	Non	Sec. 14.7.3
Classification and regression trees (CART)	Both	Discrim	Non	Sec. 16.2
Boosted model	Both	Discrim	Non	Sec. 16.4
Sparse kernelized lin/logreg (SKLR)	Both	Discrim	Non	Sec. 14.3.2
Relevance vector machine (RVM)	Both	Discrim	Non	Sec. 14.3.2
Support vector machine (SVM)	Both	Discrim	Non	Sec. 14.5
Gaussian processes (GP)	Both	Discrim	Non	Ch. 15
Smoothing splines	Regr	Discrim	Non	Section 15.4.6

Table 8.1: List of various models for classification and regression which we discuss in this book. Columns are as follows: Model name; is the model suitable for classification, regression, or both; is the model generative or discriminative; is the model parametric or non-parametric; list of sections in book which discuss the model. See also <http://pmtk3.googlecode.com/svn/trunk/docs/tutorial/html/tutSupervised.html> for the PMTK equivalents of these models. Any generative probabilistic model (e.g., HMMs, Boltzmann machines, Bayesian networks, etc.) can be turned into a classifier by using it as a class conditional density

## 8.6.2 Dealing with missing data

Sometimes some of the inputs (components of  $\mathbf{x}$ ) are not observed; this could be due to a sensor failure, or a failure to complete an entry in a survey, etc. This is called the **missing data problem** (Little. and Rubin 1987). The ability to handle missing data in a principled way is one of the biggest advantages of generative models.

To formalize our assumptions, we can associate a binary response variable  $r_i \in \{0, 1\}$  that specifies whether each value  $x_i$  is observed or not. The joint model has the form  $p(\mathbf{x}_i, r_i | \boldsymbol{\theta}, \phi) = p(r_i | \mathbf{x}_i, \phi) p(\mathbf{x}_i | \boldsymbol{\theta})$ , where  $\phi$  are the parameters controlling whether the item is observed or not.

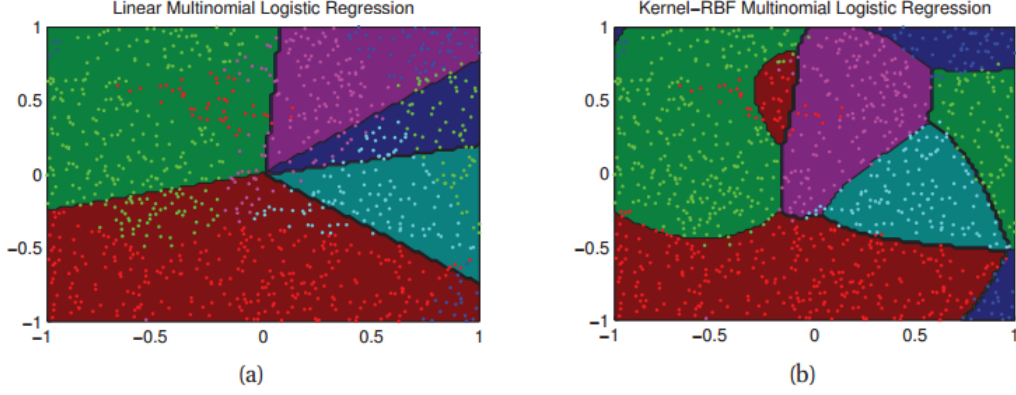


Fig. 8.1: (a) Multinomial logistic regression for 5 classes in the original feature space. (b) After basis function expansion, using RBF kernels with a bandwidth of 1, and using all the data points as centers.

- If we assume  $p(r_i|x_i, \phi) = p(r_i|\phi)$ , we say the data is **missing completely at random** or **MCAR**.
- If we assume  $p(r_i|x_i, \phi) = p(r_i|x_i^o, \phi)$ , where  $x_i^o$  is the observed part of  $x_i$ , we say the data is **missing at random** or **MAR**.
- If neither of these assumptions hold, we say the data is **not missing at random** or **NMAR**. In this case, we have to model the missing data mechanism, since the pattern of missingness is informative about the values of the missing data and the corresponding parameters. This is the case in most collaborative filtering problems, for example.

See e.g., (Marlin 2008) for further discussion. We will henceforth assume the data is MAR.

When dealing with missing data, it is helpful to distinguish the cases when there is missingness only at test time (so the training data is **complete data**), from the harder case when there is missingness also at training time. We will discuss these two cases below. Note that the class label is always missing at test time, by definition; if the class label is also sometimes missing at training time, the problem is called semi-supervised learning.

### 8.6.2.1 Missing data at test time

In a generative classifier, we can handle features that are MAR by marginalizing them out. For example, if we are missing the value of  $x_1$ , we can compute

$$p(y = c | \mathbf{x}_{2:D}, \boldsymbol{\theta}) \propto p(y = c | \boldsymbol{\theta}) p(\mathbf{x}_{2:D} | y = c, \boldsymbol{\theta}) \quad (8.23)$$

$$= \propto p(y = c | \boldsymbol{\theta}) \sum_{x_1} p(x_1, \mathbf{x}_{2:D} | y = c, \boldsymbol{\theta}) \quad (8.24)$$

Similarly, in discriminant analysis, no matter what regularization method was used to estimate the parameters, we can always analytically marginalize out the missing variables (see Section 4.3):

$$p(\mathbf{x}_{2:D} | y = c, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_{2:D} | \boldsymbol{\mu}_{c,2:D}, \boldsymbol{\Sigma}_{c,2:D}) \quad (8.25)$$

### 8.6.2.2 Missing data at training time

Missing data at training time is harder to deal with. In particular, computing the MLE or MAP estimate is no longer a simple optimization problem, for reasons discussed in Section TODO. However, soon we will study a variety of more sophisticated algorithms (such as EM algorithm, in Section 11.4) for finding approximate ML or MAP estimates in such cases.

### 8.6.3 Fishers linear discriminant analysis (FLDA) \*

TODO



## Chapter 9

# Generalized linear models and the exponential family

### 9.1 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 9.1.5).
- 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 9.1.6).
- The exponential family is at the core of generalized linear models, as discussed in Section 9.2.
- The exponential family is at the core of variational inference, as discussed in Section TODO.

#### 9.1.1 Definition

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

$$p(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp[\theta^T \phi(x)] \quad (9.1)$$

$$= h(x) \exp[\theta^T \phi(x) - A(\theta)] \quad (9.2)$$

where

$$Z(\theta) = \int h(x) \exp[\theta^T \phi(x)] dx \quad (9.3)$$

$$A(\theta) = \log Z(\theta) \quad (9.4)$$

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 (9.2) can be generalized by writing

$$p(x|\theta) = h(x) \exp[\eta(\theta)^T \phi(x) - A(\eta(\theta))] \quad (9.5)$$

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.

## 9.1.2 Examples

### 9.1.2.1 Bernoulli

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

$$\text{Ber}(x|\mu) = \mu^x(1-\mu)^{1-x} = \exp[x \log \mu + (1-x) \log(1-\mu)] \quad (9.6)$$

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

$$\begin{aligned} \text{Ber}(x|\mu) &= (1-\mu) \exp\left(x \log \frac{\mu}{1-\mu}\right) \\ \text{where } \phi(x) &= x, \theta = \log \frac{\mu}{1-\mu}, Z = \frac{1}{1-\mu} \end{aligned} \quad (9.7)$$

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

$$\mu = \text{sigm}(\theta) = \frac{1}{1 + e^{-\theta}} \quad (9.8)$$

### 9.1.2.2 Multinoulli

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

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

We can write this in exponential family form as follows:

$$\text{Cat}(\mathbf{x}|\boldsymbol{\mu}) = \exp[\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) - A(\boldsymbol{\theta})] \quad (9.9)$$

$$\boldsymbol{\theta} \triangleq \left(\log \frac{\mu_1}{\mu_K}, \dots, \log \frac{\mu_{K-1}}{\mu_K}\right) \quad (9.10)$$

$$\boldsymbol{\phi}(\mathbf{x}) \triangleq (x_1, \dots, x_{K-1}) \quad (9.11)$$

We can recover the mean parameters from the canonical parameters using

$$\mu_k = \frac{e^{\theta_k}}{1 + \sum_{j=1}^{K-1} e^{\theta_j}} \quad (9.12)$$

$$\mu_K = 1 - \frac{\sum_{j=1}^{K-1} e^{\theta_j}}{1 + \sum_{j=1}^{K-1} e^{\theta_j}} = \frac{1}{1 + \sum_{j=1}^{K-1} e^{\theta_j}} \quad (9.13)$$



and hence

$$A(\boldsymbol{\theta}) = -\log \mu_K = \log\left(1 + \sum_{j=1}^{K-1} e^{\theta_j}\right) \quad (9.14)$$

### 9.1.2.3 Univariate Gaussian

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

$$\begin{aligned} \mathcal{N}(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(x-\mu)^2\right] \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}\mu^2\right] \\ &= \frac{1}{Z(\boldsymbol{\theta})} \exp[\boldsymbol{\theta}^T \boldsymbol{\phi}(x)] \end{aligned} \quad (9.15)$$

where

$$\boldsymbol{\theta} = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right) \quad (9.16)$$

$$\boldsymbol{\phi}(x) = (x, x^2) \quad (9.17)$$

$$Z(\boldsymbol{\theta}) = \sqrt{2\pi}\sigma \exp\left(\frac{\mu^2}{2\sigma^2}\right) \quad (9.18)$$

### 9.1.2.4 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.

## 9.1.3 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>19</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>19</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  $\mathbb{E}[X^2]$ .

$$\begin{aligned}
\frac{dA}{d\theta} &= \frac{d}{d\theta} \left\{ \log \int \exp[\theta \phi(x)] h(x) dx \right\} \\
&= \frac{\frac{d}{d\theta} \int \exp[\theta \phi(x)] h(x) dx}{\int \exp[\theta \phi(x)] h(x) dx} \\
&= \frac{\int \phi(x) \exp[\theta \phi(x)] h(x) dx}{\exp(A(\theta))} \\
&= \int \phi(x) \exp[\theta \phi(x) - A(\theta)] h(x) dx \\
&= \int \phi(x) p(x) dx = \mathbb{E}[\phi(x)]
\end{aligned} \tag{9.19}$$

For the second derivative we have

$$\begin{aligned}
\frac{d^2 A}{d\theta^2} &= \int \phi(x) \exp[\theta \phi(x) - A(\theta)] h(x) [\phi(x) - A'(\theta)] dx \\
&= \int \phi(x) p(x) [\phi(x) - A'(\theta)] dx \\
&= \int \phi^2(x) p(x) dx - A'(\theta) \int \phi(x) p(x) dx \\
&= \mathbb{E}[\phi^2(x)] - \mathbb{E}[\phi(x)]^2 = \text{var}[\phi(x)]
\end{aligned} \tag{9.20}$$

In the multivariate case, we have that

$$\frac{\partial^2 A}{\partial \theta_i \partial \theta_j} = \mathbb{E}[\phi_i(x) \phi_j(x)] - \mathbb{E}[\phi_i(x)] \mathbb{E}[\phi_j(x)] \tag{9.21}$$

and hence

$$\nabla^2 A(\theta) = \text{cov}[\phi(x)] \tag{9.22}$$

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

### 9.1.4 MLE for the exponential family

The likelihood of an exponential family model has the form

$$p(\mathcal{D}|\theta) = \left[ \prod_{i=1}^N h(x_i) \right] g(\theta)^N \exp \left[ \theta^T \left( \sum_{i=1}^N \phi(x_i) \right) \right] \tag{9.23}$$

We see that the sufficient statistics are  $N$  and

$$\phi(\mathcal{D}) = \sum_{i=1}^N \phi(x_i) = \left( \sum_{i=1}^N \phi_1(x_i), \dots, \sum_{i=1}^N \phi_K(x_i) \right) \tag{9.24}$$

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.

### 9.1.5 Bayes for the exponential family

TODO

#### 9.1.5.1 Likelihood

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

## 9.2 Generalized linear models (GLMs)

Linear and logistic regression are examples of **generalized linear models**, or **GLMs** (McCullagh and Nelder 1989). These are models in which the output density is in the exponential family (Section 9.1), and in which the mean parameters are a linear combination of the inputs, passed through a possibly nonlinear function, such as the logistic function. We describe GLMs in more detail below. We focus on scalar outputs for notational simplicity. (This excludes multinomial logistic regression, but this is just to simplify the presentation.)

### 9.2.1 Basics

## 9.3 Probit regression

## 9.4 Multi-task learning



## Chapter 10

### Directed graphical models (Bayes nets)

#### 10.1 Introduction

##### 10.1.1 Chain rule

$$p(x_{1:V}) = p(x_1)p(x_2|x_1)p(x_3|x_{1:2}) \cdots p(x_N|x_{1:V-1}) \quad (10.1)$$

##### 10.1.2 Conditional independence

$X$  and  $Y$  are **conditionally independent** given  $Z$ , denoted  $X \perp Y|Z$ , iff the conditional joint can be written as a product of conditional marginals, i.e.

$$X \perp Y|Z \iff p(X, Y|Z) = p(X|Z)p(Y|Z) \quad (10.2)$$

first order **Markov assumption**: the future is independent of the past given the present,

$$x_{t+1} \perp x_{1:t-1} | x_t \quad (10.3)$$

first-order **Markov chain**

$$p(x_{1:V}) = p(x_1) \prod_{t=2}^V p(x_t | x_{t-1}) \quad (10.4)$$

##### 10.1.3 Graphical models

A **graphical model**(GM) is a way to represent a joint distribution by making CI assumptions. In particular, the nodes in the graph represent random variables, and the (lack of) edges represent CI assumptions.

There are several kinds of graphical model, depending on whether the graph is directed, undirected, or some combination of directed and undirected. In this chapter, we just study directed graphs. We consider undirected graphs in Chapter 19.

##### 10.1.4 Directed graphical model

A **directed graphical model** or **DGM** is a GM whose graph is a DAG. These are more commonly known as **Bayesian networks**. However, there is nothing inherently Bayesian about Bayesian networks: they are just a way of defining probability distributions. These models are also called **belief networks**. The term belief here refers to subjective

probability. Once again, there is nothing inherently subjective about the kinds of probability distributions represented by DGMs.

### Ordered Markov property

$$x_s \perp x_{\text{pred}(s)} \mid x_{\text{pa}(s)} \quad (10.5)$$

where  $\text{pa}(s)$  are the parents of nodes, and  $\text{pred}(s)$  are the predecessors of nodes in the DAG.

### Markov chain on a DGM

$$p(x_{1:V}|G) = \prod_{t=1}^V p(x_t | x_{\text{pa}(t)}) \quad (10.6)$$

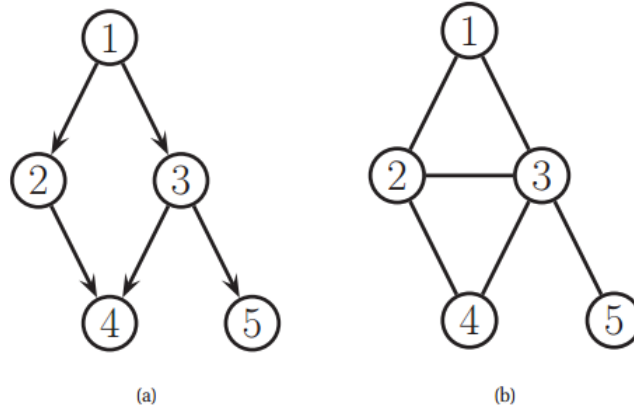


Fig. 10.1: (a) A simple DAG on 5 nodes, numbered in topological order. Node 1 is the root, nodes 4 and 5 are the leaves. (b) A simple undirected graph, with the following maximal cliques: 1,2,3, 2,3,4, 3,5.

## 10.2 Examples

### 10.2.1 Naive Bayes classifiers

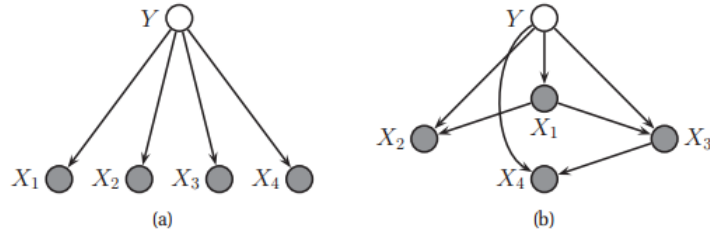


Fig. 10.2: (a) A naive Bayes classifier represented as a DGM. We assume there are  $D = 4$  features, for simplicity. Shaded nodes are observed, unshaded nodes are hidden. (b) Tree-augmented naive Bayes classifier for  $D = 4$  features. In general, the tree topology can change depending on the value of  $y$ .

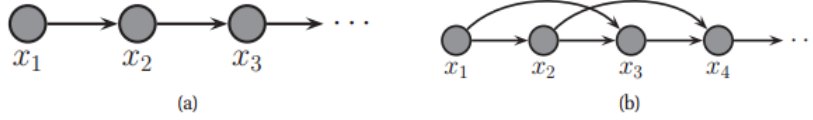


Fig. 10.3: A first and second order Markov chain.

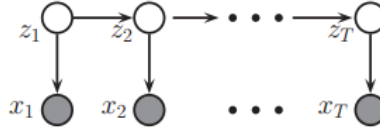


Fig. 10.4: A first-order HMM.

## 10.2.2 Markov and hidden Markov models

### 10.3 Inference

Suppose we have a set of correlated random variables with joint distribution  $p(\mathbf{x}_{1:V}|\boldsymbol{\theta})$ . Let us partition this vector into the **visible variables**  $\mathbf{x}_v$ , which are observed, and the **hidden variables**,  $\mathbf{x}_h$ , which are unobserved. Inference refers to computing the posterior distribution of the unknowns given the knowns:

$$p(\mathbf{x}_h|\mathbf{x}_v, \boldsymbol{\theta}) = \frac{p(\mathbf{x}_h, \mathbf{x}_v|\boldsymbol{\theta})}{p(\mathbf{x}_v|\boldsymbol{\theta})} = \frac{p(\mathbf{x}_h, \mathbf{x}_v|\boldsymbol{\theta})}{\sum_{\mathbf{x}'_h} p(\mathbf{x}'_h, \mathbf{x}_v|\boldsymbol{\theta})} \quad (10.7)$$

Sometimes only some of the hidden variables are of interest to us. So let us partition the hidden variables into **query variables**,  $\mathbf{x}_q$ , whose value we wish to know, and the remaining **nuisance variables**,  $\mathbf{x}_n$ , which we are not interested in. We can compute what we are interested in by **marginalizing out** the nuisance variables:

$$p(\mathbf{x}_q|\mathbf{x}_v, \boldsymbol{\theta}) = \sum_{\mathbf{x}_n} p(\mathbf{x}_q, \mathbf{x}_n|\mathbf{x}_v, \boldsymbol{\theta}) \quad (10.8)$$

### 10.4 Learning

MAP estimate:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p(\mathbf{x}_{i,v}|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \quad (10.9)$$

#### 10.4.1 Learning from complete data

If all the variables are fully observed in each case, so there is no missing data and there are no hidden variables, we say the data is **complete**. For a DGM with complete data, the likelihood is given by

$$p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{i=1}^N p(\mathbf{x}_i|\boldsymbol{\theta}) = \prod_{i=1}^N \prod_{t=1}^V p(\mathbf{x}_{it}|\mathbf{x}_{i,\text{pa}(t)}, \boldsymbol{\theta}_t) = \prod_{t=1}^V p(\mathcal{D}_t|\boldsymbol{\theta}_t) \quad (10.10)$$

where  $\mathcal{D}_t$  is the data associated with node  $t$  and its parents, i.e., the  $t$ 'th family.

Now suppose that the prior factorizes as well:

$$p(\boldsymbol{\theta}) = \prod_{t=1}^V p(\boldsymbol{\theta}_t) \quad (10.11)$$

Then clearly the posterior also factorizes:

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}) = \prod_{t=1}^V p(\mathcal{D}_t|\boldsymbol{\theta}_t)p(\boldsymbol{\theta}_t) \quad (10.12)$$

## 10.4.2 Learning with missing and/or latent variables

If we have missing data and/or hidden variables, the likelihood no longer factorizes, and indeed it is no longer convex, as we explain in detail in Section TODO. This means we will usually can only compute a locally optimal ML or MAP estimate. Bayesian inference of the parameters is even harder. We discuss suitable approximate inference techniques in later chapters.

## 10.5 Conditional independence properties of DGMs

### 10.5.1 d-separation and the Bayes Ball algorithm (global Markov properties)

1. P contains a chain

$$p(x, z|y) = \frac{p(x, y, z)}{p(y)} = \frac{p(x)p(y|x)p(z|y)}{p(y)} = \frac{p(x, y)p(z|y)}{p(y)} = p(x|y)p(z|y) \quad (10.13)$$

2. P contains a fork

$$p(x, z|y) = \frac{p(x, y, z)}{p(y)} = \frac{p(y)p(x|y)p(z|y)}{p(y)} = p(x|y)p(z|y) \quad (10.14)$$

3. P contains v-structure

$$p(x, z|y) = \frac{p(x, y, z)}{p(y)} = \frac{p(x)p(z)p(y|x, z)}{p(y)} \neq p(x|y)p(z|y) \quad (10.15)$$

### 10.5.2 Other Markov properties of DGMs

### 10.5.3 Markov blanket and full conditionals

$$mb(t) = ch(t) \cup pa(t) \cup cop_a(t) \quad (10.16)$$

### 10.5.4 Multinoulli Learning

Multinoulli Distribution

$$Cat(x|\mu) = \prod_{k=1}^K \mu_k^{x_k} \quad (10.17)$$

then from ?? and 10.17:

$$p(x|G, \boldsymbol{\theta}) = \prod_{v=1}^V \prod_{c=1}^{C_v} \prod_{k=1}^K \theta_{vck}^{y_{vck}} \quad (10.18)$$



Likelihood

$$p(D|G, \theta) = \prod_{n=1}^N p(x_n|G, \theta) = \prod_{n=1}^N \prod_{v=1}^V \prod_{c=1}^{C_{nv}} \prod_{k=1}^K \theta_{vck}^{y_{nvck}} \quad (10.19)$$

where  $y_{nv} = f(pa(x_{nv}))$ ,  $f(x)$  is a map from  $x$  to a vector, there is only one element in the vector is 1.

## 10.6 Influence (decision) diagrams \*



# Chapter 11

## Mixture models and the EM algorithm

### 11.1 Latent variable models

In Chapter 10 we showed how graphical models can be used to define high-dimensional joint probability distributions. The basic idea is to model dependence between two variables by adding an edge between them in the graph. (Technically the graph represents conditional independence, but you get the point.)

An alternative approach is to assume that the observed variables are correlated because they arise from a hidden common cause. Model with hidden variables are also known as **latent variable models** or **LVMs**. As we will see in this chapter, such models are harder to fit than models with no latent variables. However, they can have significant advantages, for two main reasons.

- First, LVMs often have fewer parameters than models that directly represent correlation in the visible space.
- Second, the hidden variables in an LVM can serve as a **bottleneck**, which computes a compressed representation of the data. This forms the basis of unsupervised learning, as we will see. Figure 11.1 illustrates some generic LVM structures that can be used for this purpose.

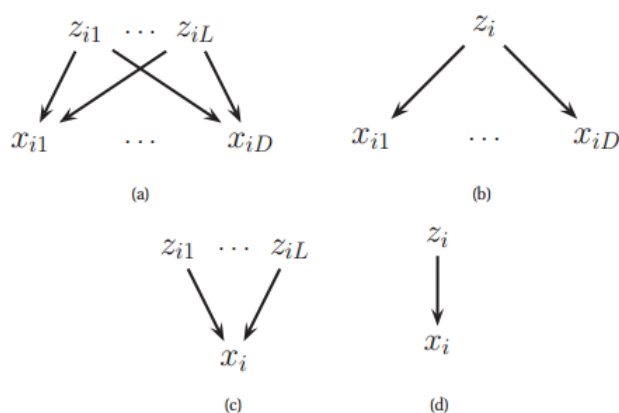


Fig. 11.1: A latent variable model represented as a DGM. (a) Many-to-many. (b) One-to-many. (c) Many-to-one. (d) One-to-one.

### 11.2 Mixture models

The simplest form of LVM is when  $z_i \in \{1, \dots, K\}$ , representing a discrete latent state. We will use a discrete prior for this,  $p(z_i) = \text{Cat}(\pi)$ . For the likelihood, we use  $p(x_i|z_i = k) = p_k(x_i)$ , where  $p_k$  is the  $k$ 'th **base distribution** for the

observations; this can be of any type. The overall model is known as a **mixture model**, since we are mixing together the  $K$  base distributions as follows:

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_i|\boldsymbol{\theta}) \quad (11.1)$$

Depending on the form of the likelihood  $p(\mathbf{x}_i|\mathbf{z}_i)$  and the prior  $p(\mathbf{z}_i)$ , we can generate a variety of different models, as summarized in Table 11.1.

$p(\mathbf{x}_i \mathbf{z}_i)$	$p(\mathbf{z}_i)$	Name	Section
MVN	Discrete	Mixture of Gaussians	11.2.1
Prod. Discrete	Discrete	Mixture of multinomials	11.2.2
Prod. Gaussian	Prod. Gaussian	Factor analysis/ probabilistic PCA	12.1.5
Prod. Gaussian	Prod. Laplace	Probabilistic ICA/ sparse coding	12.6
Prod. Discrete	Prod. Gaussian	Multinomial PCA	27.2.3
Prod. Discrete	Dirichlet	Latent Dirichlet allocation	27.3
Prod. Noisy-OR	Prod. Bernoulli	BN20/ QMR	10.2.3
Prod. Bernoulli	Prod. Bernoulli	Sigmoid belief net	27.7

Table 11.1: Summary of some popular directed latent variable models. Here Prod means product, so Prod. Discrete in the likelihood means a factored distribution of the form  $\prod_j \text{Cat}(x_{ij}|\mathbf{z}_i)$ , and Prod. Gaussian means a factored distribution of the form  $\prod_j \mathcal{N}(x_{ij}|\mathbf{z}_i)$ .

### 11.2.1 Mixtures of Gaussians

$$p_k(\mathbf{x}_i|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (11.2)$$

### 11.2.2 Mixtures of multinoullis

$$p_k(\mathbf{x}_i|\boldsymbol{\theta}) = \prod_{j=1}^D \text{Ber}(x_{ij}|\mu_{jk}) = \prod_{j=1}^D \mu_{jk}^{x_{ij}} (1 - \mu_{jk})^{1-x_{ij}} \quad (11.3)$$

where  $\mu_{jk}$  is the probability that bit  $j$  turns on in cluster  $k$ .

The latent variables do not have to any meaning, we might simply introduce latent variables in order to make the model more powerful. For example, one can show that the mean and covariance of the mixture distribution are given by

$$\mathbb{E}[\mathbf{x}] = \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k \quad (11.4)$$

$$\text{Cov}[\mathbf{x}] = \sum_{k=1}^K \pi_k (\boldsymbol{\Sigma}_k + \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T) - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T \quad (11.5)$$

where  $\boldsymbol{\Sigma}_k = \text{diag}(\mu_{jk}(1 - \mu_{jk}))$ . So although the component distributions are factorized, the joint distribution is not. Thus the mixture distribution can capture correlations between variables, unlike a single product-of-Bernoullis model.

### 11.2.3 Using mixture models for clustering

There are two main applications of mixture models, black-box density model(see Section 14.7.3 TODO) and clustering(see Chapter 25 TODO).

#### Soft clustering

$$r_{ik} \triangleq p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}) = \frac{p(z_i = k, \mathbf{x}_i | \boldsymbol{\theta})}{p(\mathbf{x}_i | \boldsymbol{\theta})} = \frac{p(z_i = k | \boldsymbol{\theta}) p(\mathbf{x}_i | z_i = k, \boldsymbol{\theta})}{\sum_{k'=1}^K p(z_i = k' | \boldsymbol{\theta}) p(\mathbf{x}_i | z_i = k', \boldsymbol{\theta})} \quad (11.6)$$

where  $r_{ik}$  is known as the **responsibility** of cluster  $k$  for point  $i$ .

#### Hard clustering

$$z_i^* \triangleq \arg \max_k r_{ik} = \arg \max_k p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}) \quad (11.7)$$

The difference between generative classifiers and mixture models only arises at training time: in the mixture case, we never observe  $z_i$ , whereas with a generative classifier, we do observe  $y_i$ (which plays the role of  $z_i$ ).

### 11.2.4 Mixtures of experts

Section 14.7.3 TODO described how to use mixture models in the context of generative classifiers. We can also use them to create discriminative models for classification and regression. For example, consider the data in Figure 11.2(a). It seems like a good model would be three different linear regression functions, each applying to a different part of the input space. We can model this by allowing the mixing weights and the mixture densities to be input-dependent:

$$p(y_i | \mathbf{x}_i, z_i = k, \boldsymbol{\theta}) = \mathcal{N}(y_i | \mathbf{w}_k^T \mathbf{x}, \sigma_k^2) \quad (11.8)$$

$$p(z_i | \mathbf{x}_i, \boldsymbol{\theta}) = \text{Cat}(z_i | \mathcal{S}(\mathbf{V}^T \mathbf{x}_i)) \quad (11.9)$$

See Figure 11.3(a) for the DGM.

This model is called a **mixture of experts** or MoE (Jordan and Jacobs 1994). The idea is that each submodel is considered to be an expert in a certain region of input space. The function  $p(z_i | \mathbf{x}_i, \boldsymbol{\theta})$  is called a **gating function**, and decides which expert to use, depending on the input values. For example, Figure 11.2(b) shows how the three experts have carved up the 1d input space, Figure 11.2(a) shows the predictions of each expert individually (in this case, the experts are just linear regression models), and Figure 11.2(c) shows the overall prediction of the model, obtained using

$$p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = \sum_{k=1}^K p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}) p(y_i | \mathbf{x}_i, z_i = k, \boldsymbol{\theta}) \quad (11.10)$$

We discuss how to fit this model in Section TODO 11.4.3.

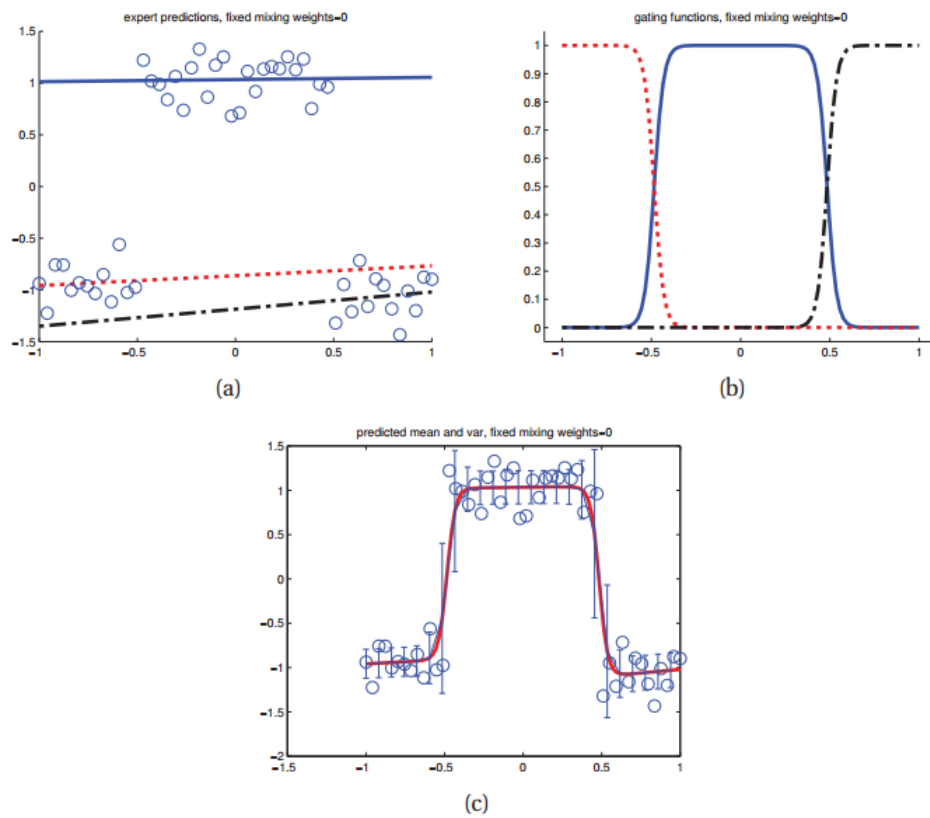


Fig. 11.2: (a) Some data fit with three separate regression lines. (b) Gating functions for three different experts. (c) The conditionally weighted average of the three expert predictions.

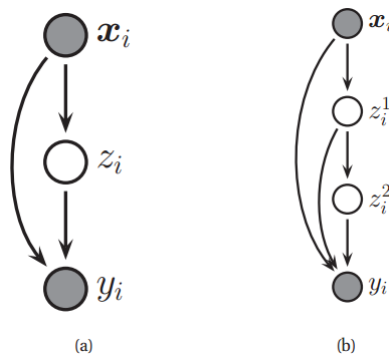


Fig. 11.3: (a) A mixture of experts. (b) A hierarchical mixture of experts.

## 11.3 Parameter estimation for mixture models

### 11.3.1 Unidentifiability

### 11.3.2 Computing a MAP estimate is non-convex

## 11.4 The EM algorithm

### 11.4.1 Introduction

For many models in machine learning and statistics, computing the ML or MAP parameter estimate is easy provided we observe all the values of all the relevant random variables, i.e., if we have complete data. However, if we have missing data and/or latent variables, then computing the ML/MAP estimate becomes hard.

One approach is to use a generic gradient-based optimizer to find a local minimum of the  $\text{NLL}(\theta)$ . However, we often have to enforce constraints, such as the fact that covariance matrices must be positive definite, mixing weights must sum to one, etc., which can be tricky. In such cases, it is often much simpler (but not always faster) to use an algorithm called **expectation maximization**, or **EM** for short (Dempster et al. 1977; Meng and van Dyk 1997; McLachlan and Krishnan 1997). This is an efficient iterative algorithm to compute the ML or MAP estimate in the presence of missing or hidden data, often with closed-form updates at each step. Furthermore, the algorithm automatically enforces the required constraints.

See Table 11.2 for a summary of the applications of EM in this book.

Model	Section
Mix. Gaussians	11.4.2
Mix. experts	11.4.3
Factor analysis	12.1.5
Student T	11.4.5
Probit regression	11.4.6
DGM with hidden variables	11.4.4
MVN with missing data	11.6.1
HMMs	17.5.2
Shrinkage estimates of Gaussian means	Exercise 11.13

Table 11.2: Some models discussed in this book for which EM can be easily applied to find the ML/ MAP parameter estimate.

### 11.4.2 Basic idea

EM exploits the fact that if the data were fully observed, then the ML/ MAP estimate would be easy to compute. In particular, each iteration of the EM algorithm consists of two processes: The E-step, and the M-step.

- In the **E-step**, the missing data are inferred given the observed data and current estimate of the model parameters. This is achieved using the conditional expectation, explaining the choice of terminology.
- In the **M-step**, the likelihood function is maximized under the assumption that the missing data are known. The missing data inferred from the E-step are used in lieu of the actual missing data.

Let  $x_i$  be the visible or observed variables in case  $i$ , and let  $z_i$  be the hidden or missing variables. The goal is to maximize the log likelihood of the observed data:

$$\ell(\theta) = \log p(\mathcal{D}|\theta) = \sum_{i=1}^N \log p(x_i|\theta) = \sum_{i=1}^N \log \sum_{z_i} p(x_i, z_i|\theta) \quad (11.11)$$

Unfortunately this is hard to optimize, since the log cannot be pushed inside the sum. EM gets around this problem as follows. Define the **complete data log likelihood** to be

$$\ell_c(\theta) = \sum_{i=1}^N \log p(x_i, z_i|\theta) \quad (11.12)$$

This cannot be computed, since  $z_i$  is unknown. So let us define the **expected complete data log likelihood** as follows:

$$Q(\theta, \theta^{t-1}) \triangleq \mathbb{E}_{z|\mathcal{D}, \theta^{t-1}} [\ell_c(\theta)] = \mathbb{E} [\ell_c(\theta) | \mathcal{D}, \theta^{t-1}] \quad (11.13)$$

where  $t$  is the current iteration number.  $Q$  is called the **auxiliary function** (see Section 11.4.8 for derivation). The expectation is taken wrt the old parameters,  $\theta^{t-1}$ , and the observed data  $\mathcal{D}$ . The goal of the E-step is to compute  $Q(\theta, \theta^{t-1})$ , or rather, the parameters inside of it which the MLE (or MAP) depends on; these are known as the **expected sufficient statistics** or **ESS**. In the M-step, we optimize the  $Q$  function wrt  $\theta$ :

$$\theta^t = \arg \max_{\theta} Q(\theta, \theta^{t-1}) \quad (11.14)$$

To perform MAP estimation, we modify the M-step as follows:

$$\theta^t = \arg \max_{\theta} Q(\theta, \theta^{t-1}) + \log p(\theta) \quad (11.15)$$

The E step remains unchanged.

In summary, the EM algorithm's pseudo code is as follows

```

input : observed data  $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ , joint distribution  $P(x, z|\theta)$ 
output: model's parameters  $\theta$ 
// 1. identify hidden variables  $z$ , write out the log likelihood function  $\ell(x, z|\theta)$ 
 $\theta^{(0)} = \dots$  // initialize
while (!convergence) do
    // 2. E-step: plug in  $P(x, z|\theta)$ , derive the formula of  $Q(\theta, \theta^{t-1})$ 
     $Q(\theta, \theta^{t-1}) = \mathbb{E} [\ell_c(\theta) | \mathcal{D}, \theta^{t-1}]$ 
    // 3. M-step: find  $\theta$  that maximizes the value of  $Q(\theta, \theta^{t-1})$ 
     $\theta^t = \arg \max_{\theta} Q(\theta, \theta^{t-1})$ 
end
```

**Algorithm 1:** EM algorithm

Below we explain how to perform the E and M steps for several simple models, that should make things clearer.

### 11.4.3 EM for GMMs

#### 11.4.3.1 Auxiliary function



$$\begin{aligned}
Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1}) &= \mathbb{E}_{z|\mathcal{D}, \boldsymbol{\theta}^{t-1}} [\ell_c(\boldsymbol{\theta})] \\
&= \mathbb{E}_{z|\mathcal{D}, \boldsymbol{\theta}^{t-1}} \left[ \sum_{i=1}^N \log p(\mathbf{x}_i, z_i | \boldsymbol{\theta}) \right] \\
&= \sum_{i=1}^N \mathbb{E}_{z|\mathcal{D}, \boldsymbol{\theta}^{t-1}} \left\{ \log \left[ \prod_{k=1}^K (\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k))^{\mathbb{I}(z_i=k)} \right] \right\} \\
&= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E}[\mathbb{I}(z_i = k)] \log [\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k)] \\
&= \sum_{i=1}^N \sum_{k=1}^K p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{t-1}) \log [\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k)] \\
&= \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log p(\mathbf{x}_i | \boldsymbol{\theta}_k)
\end{aligned} \tag{11.16}$$

where  $r_{ik} \triangleq \mathbb{E}[\mathbb{I}(z_i = k)] = p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{t-1})$  is the **responsibility** that cluster  $k$  takes for data point  $i$ . This is computed in the E-step, described below.

#### 11.4.3.2 E-step

The E-step has the following simple form, which is the same for any mixture model:

$$r_{ik} = p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{t-1}) = \frac{p(z_i = k, \mathbf{x}_i | \boldsymbol{\theta}^{t-1})}{p(\mathbf{x}_i | \boldsymbol{\theta}^{t-1})} = \frac{\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k^{t-1})}{\sum_{k'=1}^K \pi_{k'} p(\mathbf{x}_i | \boldsymbol{\theta}_{k'}^{t-1})} \tag{11.17}$$

#### 11.4.3.3 M-step

In the M-step, we optimize  $Q$  wrt  $\boldsymbol{\pi}$  and  $\boldsymbol{\theta}_k$ .

For  $\boldsymbol{\pi}$ , grouping together only the terms that depend on  $\pi_k$ , we find that we need to maximize  $\sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \pi_k$ .

However, there is an additional constraint  $\sum_{k=1}^K \pi_k = 1$ , since they represent the probabilities  $\pi_k = P(z_i = k)$ . To deal with the constraint we construct the Lagrangian

$$\mathcal{L}(\boldsymbol{\pi}) = \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \pi_k + \beta \left( \sum_{k=1}^K \pi_k - 1 \right)$$

where  $\beta$  is the Lagrange multiplier. Taking derivatives, we find

$$\hat{\pi}_k = \frac{\sum_{i=1}^N \hat{r}_{ik}}{N} \tag{11.18}$$

This is the same for any mixture model, whereas  $\boldsymbol{\theta}_k$  depends on the form of  $p(\mathbf{x} | \boldsymbol{\theta}_k)$ .

For  $\boldsymbol{\theta}_k$ , plug in the pdf to Equation (11.16)

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1}) = \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \pi_k - \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^K r_{ik} [\log |\boldsymbol{\Sigma}_k| + (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)]$$

Take partial derivatives of  $Q$  wrt  $\boldsymbol{\mu}_k$ ,  $\boldsymbol{\Sigma}_k$  and let them equal to 0, we can get

$$\begin{aligned}
\frac{\partial Q}{\partial \boldsymbol{\mu}_k} &= -\frac{1}{2} \sum_{i=1}^N r_{ik} [(\boldsymbol{\Sigma}_k^{-1} + \boldsymbol{\Sigma}_k^{-T})(\mathbf{x}_i - \boldsymbol{\mu}_k)] \\
&= -\sum_{i=1}^N r_{ik} [\boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k)] = 0 \Rightarrow \\
\hat{\boldsymbol{\mu}}_k &= \frac{\sum_{i=1}^N \hat{r}_{ik} \mathbf{x}_i}{\sum_{i=1}^N \hat{r}_{ik}}
\end{aligned} \tag{11.19}$$

$$\begin{aligned}
\frac{\partial Q}{\partial \boldsymbol{\Sigma}_k} &= -\frac{1}{2} \sum_{i=1}^N r_{ik} \left[ \frac{1}{\boldsymbol{\Sigma}_k} - \frac{1}{\boldsymbol{\Sigma}_k^2} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \right] = 0 \Rightarrow \\
\hat{\boldsymbol{\Sigma}}_k &= \frac{\sum_{i=1}^N \hat{r}_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T}{\sum_{i=1}^N \hat{r}_{ik}} = \frac{\sum_{i=1}^N \hat{r}_{ik} \mathbf{x}_i \mathbf{x}_i^T}{\sum_{i=1}^N \hat{r}_{ik}} - \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T
\end{aligned} \tag{11.20}$$

#### 11.4.3.4 Algorithm pseudo code

```

input : observed data  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ , GMM
output: GMM's parameters  $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ 
// 1. initialize
 $\boldsymbol{\pi}^{(0)} = \dots$ 
 $\boldsymbol{\mu}^{(0)} = \dots$ 
 $\boldsymbol{\Sigma}^{(0)} = \dots$ 
t = 0
while (!convergency) do
  // 2. E-step
   $\hat{r}_{ik} = \frac{\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k^{t-1})}{\sum_{k'=1}^K \pi_{k'} p(\mathbf{x}_i | \boldsymbol{\mu}_{k'}^{t-1}, \boldsymbol{\Sigma}_{k'}^{t-1})}$ 
  // 3. M-step
   $\hat{\boldsymbol{\pi}}_k = \frac{\sum_{i=1}^N \hat{r}_{ik}}{N}$ 
   $\hat{\boldsymbol{\mu}}_k = \frac{\sum_{i=1}^N \hat{r}_{ik} \mathbf{x}_i}{\sum_{i=1}^N \hat{r}_{ik}}$ 
   $\hat{\boldsymbol{\Sigma}}_k = \frac{\sum_{i=1}^N \hat{r}_{ik} \mathbf{x}_i \mathbf{x}_i^T}{\sum_{i=1}^N \hat{r}_{ik}} - \hat{\boldsymbol{\mu}}_k \hat{\boldsymbol{\mu}}_k^T$ 
  ++t
end

```

**Algorithm 2:** EM algorithm for GMM

#### 11.4.3.5 MAP estimation

As usual, the MLE may overfit. The overfitting problem is particularly severe in the case of GMMs. An easy solution to this is to perform MAP estimation. The new auxiliary function is the expected complete data log-likelihood plus the log prior:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1}) = \left[ \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log p(\mathbf{x}_i | \boldsymbol{\theta}_k) \right] + \log p(\boldsymbol{\pi}) + \sum_{k=1}^K \log p(\boldsymbol{\theta}_k) \tag{11.21}$$

It is natural to use conjugate priors.

$$\begin{aligned}
p(\boldsymbol{\pi}) &= \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}) \\
p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) &= \text{NIW}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k | \mathbf{m}_0, \boldsymbol{\kappa}_0, \mathbf{v}_0, \mathbf{S}_0)
\end{aligned}$$

From Equation (3.29) and Section 4.4.3, the MAP estimate is given by

$$\hat{\pi}_k = \frac{\sum_{i=1}^N r_{ik} + \alpha_k - 1}{N + \sum_{k=1}^K \alpha_k - K} \quad (11.22)$$

$$\hat{\mu}_k = \frac{\sum_{i=1}^N r_{ik} \mathbf{x}_i + \kappa_0 \mathbf{m}_0}{\sum_{i=1}^N r_{ik} + \kappa_0} \quad (11.23)$$

$$\hat{\Sigma}_k = \frac{\mathbf{S}_0 + \mathbf{S}_k + \frac{\kappa_0 r_k}{\kappa_0 + r_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T}{v_0 + r_k + D + 2} \quad (11.24)$$

$$\text{where } r_k \triangleq \sum_{i=1}^N r_{ik}, \bar{\mathbf{x}}_k \triangleq \frac{\sum_{i=1}^N r_{ik} \mathbf{x}_i}{r_k}, \mathbf{S}_k \triangleq \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \bar{\mathbf{x}}_k)(\mathbf{x}_i - \bar{\mathbf{x}}_k)^T$$

#### 11.4.3.6 K-means algorithm

#### 11.4.4 EM for mixture of experts

#### 11.4.5 EM for DGMs with hidden variables

#### 11.4.6 EM for the Student distribution \*

#### 11.4.7 EM for probit regression \*

#### 11.4.8 Derivation of the $Q$ function

**Theorem 11.1. (Jensen's inequality)** Let  $f$  be a convex function (see Section [A.1](#)) defined on a convex set  $\mathcal{S}$ . If  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathcal{S}$  and  $\lambda_1, \lambda_2, \dots, \lambda_n \geq 0$  with  $\sum_{i=1}^n \lambda_i = 1$ ,

$$f\left(\sum_{i=1}^n \lambda_i \mathbf{x}_i\right) \leq \sum_{i=1}^n \lambda_i f(\mathbf{x}_i) \quad (11.25)$$

**Proposition 11.1.**

$$\log\left(\sum_{i=1}^n \lambda_i \mathbf{x}_i\right) \geq \sum_{i=1}^n \lambda_i \log(\mathbf{x}_i) \quad (11.26)$$

Now let's proof why the  $Q$  function should look like Equation ([11.13](#)):

$$\begin{aligned}
\ell(\boldsymbol{\theta}) &= \log P(\mathcal{D}|\boldsymbol{\theta}) \\
&= \log \sum_{\mathbf{z}} P(\mathcal{D}, \mathbf{z}|\boldsymbol{\theta}) \\
&= \log \sum_{\mathbf{z}} P(\mathcal{D}|\mathbf{z}, \boldsymbol{\theta}) P(\mathbf{z}|\boldsymbol{\theta}) \\
\ell(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}^{t-1}) &= \log \left[ \sum_{\mathbf{z}} P(\mathcal{D}|\mathbf{z}, \boldsymbol{\theta}) P(\mathbf{z}|\boldsymbol{\theta}) \right] - \log P(\mathcal{D}|\boldsymbol{\theta}^{t-1}) \\
&= \log \left[ \sum_{\mathbf{z}} P(\mathcal{D}|\mathbf{z}, \boldsymbol{\theta}) P(\mathbf{z}|\boldsymbol{\theta}) \frac{P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1})}{P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1})} \right] - \log P(\mathcal{D}|\boldsymbol{\theta}^{t-1}) \\
&= \log \left[ \sum_{\mathbf{z}} P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) \frac{P(\mathcal{D}|\mathbf{z}, \boldsymbol{\theta}) P(\mathbf{z}|\boldsymbol{\theta})}{P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1})} \right] - \log P(\mathcal{D}|\boldsymbol{\theta}^{t-1}) \\
&\geq \sum_{\mathbf{z}} P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) \log \left[ \frac{P(\mathcal{D}|\mathbf{z}, \boldsymbol{\theta}) P(\mathbf{z}|\boldsymbol{\theta})}{P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1})} \right] - \log P(\mathcal{D}|\boldsymbol{\theta}^{t-1}) \\
&= \sum_{\mathbf{z}} P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) \log \left[ \frac{P(\mathcal{D}|\mathbf{z}, \boldsymbol{\theta}) P(\mathbf{z}|\boldsymbol{\theta})}{P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) P(\mathcal{D}|\boldsymbol{\theta}^{t-1})} \right] \\
B(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1}) \ell(\boldsymbol{\theta}^{t-1}) &\triangleq \sum_{\mathbf{z}} P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) \log \left[ \frac{P(\mathcal{D}|\mathbf{z}, \boldsymbol{\theta}) P(\mathbf{z}|\boldsymbol{\theta})}{P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) P(\mathcal{D}|\boldsymbol{\theta}^{t-1})} \right] \\
&\Rightarrow
\end{aligned} \tag{11.27}$$

$$\begin{aligned}
\boldsymbol{\theta}^t &= \arg \max_{\boldsymbol{\theta}} B(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1}) \\
&= \arg \max_{\boldsymbol{\theta}} \left\{ \ell(\boldsymbol{\theta}^{t-1}) + \sum_{\mathbf{z}} P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) \log \left[ \frac{P(\mathcal{D}|\mathbf{z}, \boldsymbol{\theta}) P(\mathbf{z}|\boldsymbol{\theta})}{P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) P(\mathcal{D}|\boldsymbol{\theta}^{t-1})} \right] \right\} \\
&\text{Now drop terms which are constant w.r.t. } \boldsymbol{\theta} \\
&= \arg \max_{\boldsymbol{\theta}} \left\{ \sum_{\mathbf{z}} P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) \log [P(\mathcal{D}|\mathbf{z}, \boldsymbol{\theta}) P(\mathbf{z}|\boldsymbol{\theta})] \right\} \\
&= \arg \max_{\boldsymbol{\theta}} \left\{ \sum_{\mathbf{z}} P(\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}) \log [P(\mathcal{D}, \mathbf{z}|\boldsymbol{\theta})] \right\} \\
&= \arg \max_{\boldsymbol{\theta}} \left\{ \mathbb{E}_{\mathbf{z}|\mathcal{D}, \boldsymbol{\theta}^{t-1}} \log [P(\mathcal{D}, \mathbf{z}|\boldsymbol{\theta})] \right\} \\
&\triangleq \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1})
\end{aligned} \tag{11.28}$$

$$\tag{11.29}$$

## 11.4.9 Convergence of the EM Algorithm \*

### 11.4.9.1 Expected complete data log likelihood is a lower bound

Note that  $\ell(\boldsymbol{\theta}) \geq B(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1})$ , and  $\ell(\boldsymbol{\theta}^{t-1}) \geq B(\boldsymbol{\theta}^{t-1}, \boldsymbol{\theta}^{t-1})$ , which means  $B(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1})$  is an lower bound of  $\ell(\boldsymbol{\theta})$ . If we maximize  $B(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1})$ , then  $\ell(\boldsymbol{\theta})$  gets maximized, see Figure 11.4.

Since the expected complete data log likelihood  $Q$  is derived from  $B(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1})$  by dropping terms which are constant w.r.t.  $\boldsymbol{\theta}$ , so it is also a lower bound to a lower bound of  $\ell(\boldsymbol{\theta})$ .

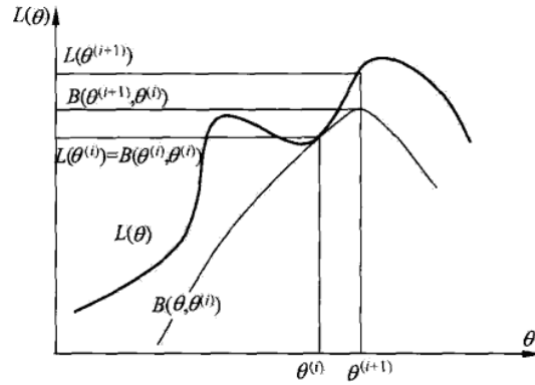


Fig. 11.4: Graphical interpretation of a single iteration of the EM algorithm: The function  $B(\theta, \theta^{t-1})$  is bounded above by the log likelihood function  $\ell(\theta)$ . The functions are equal at  $\theta = \theta^{t-1}$ . The EM algorithm chooses  $\theta^{t-1}$  as the value of  $\theta$  for which  $B(\theta, \theta^{t-1})$  is a maximum. Since  $\ell(\theta) \geq B(\theta, \theta^{t-1})$  increasing  $B(\theta, \theta^{t-1})$  ensures that the value of the log likelihood function  $\ell(\theta)$  is increased at each step.

#### 11.4.9.2 EM monotonically increases the observed data log likelihood

#### 11.4.10 Generalization of EM Algorithm \*

EM algorithm can be interpreted as F function's maximization-maximization algorithm, based on this interpretation there are many variations and generalization, e.g., generalized EM Algorithm(GEM).

##### 11.4.10.1 F function's maximization-maximization algorithm

**Definition 11.1.** Given the probability distribution of the hidden variable  $Z$  is  $\tilde{P}(Z)$ , define **F function** as the following:

$$F(\tilde{P}, \theta) = \mathbb{E}_{\tilde{P}}[\log P(X, Z|\theta)] + H(\tilde{P}) \quad (11.30)$$

Where  $H(\tilde{P}) = -\mathbb{E}_{\tilde{P}} \log \tilde{P}(Z)$ , which is  $\tilde{P}(Z)$ 's entropy. Usually we assume that  $P(X, Z|\theta)$  is continuous w.r.t.  $\theta$ , therefore  $F(\tilde{P}, \theta)$  is continuous w.r.t.  $\tilde{P}$  and  $\theta$ .

**Lemma 11.1.** For a fixed  $\theta$ , there is only one distribution  $\tilde{P}_\theta$  which maximizes  $F(\tilde{P}, \theta)$

$$\tilde{P}_\theta(Z) = P(Z|X, \theta) \quad (11.31)$$

and  $\tilde{P}_\theta$  is continuous w.r.t.  $\theta$ .

*Proof.* Given a fixed  $\theta$ , we can get  $\tilde{P}_\theta$  which maximizes  $F(\tilde{P}, \theta)$ . we construct the Lagrangian

$$\mathcal{L}(\tilde{P}, \theta) = \mathbb{E}_{\tilde{P}}[\log P(X, Z|\theta)] - \mathbb{E}_{\tilde{P}} \log \tilde{P}_\theta(Z) + \lambda \left[ 1 - \sum_Z \tilde{P}(Z) \right] \quad (11.32)$$

Take partial derivative with respect to  $\tilde{P}_\theta(Z)$  then we get

$$\frac{\partial \mathcal{L}}{\partial \tilde{P}_\theta(Z)} = \log P(X, Z|\theta) - \log \tilde{P}_\theta(Z) - 1 - \lambda$$

Let it equal to 0, we can get

$$\lambda = \log P(X, Z|\theta) - \log \tilde{P}_\theta(Z) - 1$$

Then we can derive that  $\tilde{P}_\theta(Z)$  is proportional to  $P(X, Z|\theta)$

$$\begin{aligned}
\frac{P(X, Z|\boldsymbol{\theta})}{\tilde{P}_\theta(Z)} &= e^{1+\lambda} \\
\Rightarrow \tilde{P}_\theta(Z) &= \frac{P(X, Z|\boldsymbol{\theta})}{e^{1+\lambda}} \\
\sum_Z \tilde{P}_\theta(Z) = 1 &\Rightarrow \sum_Z \frac{P(X, Z|\boldsymbol{\theta})}{e^{1+\lambda}} = 1 \Rightarrow P(X|\boldsymbol{\theta}) = e^{1+\lambda} \\
\tilde{P}_\theta(Z) &= \frac{P(X, Z|\boldsymbol{\theta})}{e^{1+\lambda}} = \frac{P(X, Z|\boldsymbol{\theta})}{P(X|\boldsymbol{\theta})} = P(Z|X, \boldsymbol{\theta})
\end{aligned}$$

**Lemma 11.2.** If  $\tilde{P}_\theta(Z) = P(Z|X, \boldsymbol{\theta})$ , then

$$F(\tilde{P}, \boldsymbol{\theta}) = \log P(X|\boldsymbol{\theta}) \quad (11.33)$$

**Theorem 11.2.** One iteration of EM algorithm can be implemented as  $F$  function's maximization-maximization.

Assume  $\boldsymbol{\theta}^{t-1}$  is the estimation of  $\boldsymbol{\theta}$  in the  $(t-1)$ -th iteration,  $\tilde{P}^{t-1}$  is the estimation of  $\tilde{P}$  in the  $(t-1)$ -th iteration. Then in the  $t$ -th iteration two steps are:

1. for fixed  $\boldsymbol{\theta}^{t-1}$ , find  $\tilde{P}^t$  that maximizes  $F(\tilde{P}, \boldsymbol{\theta}^{t-1})$ ;
2. for fixed  $\tilde{P}^t$ , find  $\boldsymbol{\theta}^t$  that maximizes  $F(\tilde{P}^t, \boldsymbol{\theta})$ .

*Proof.* (1) According to Lemma 11.1, we can get

$$\tilde{P}^t(Z) = P(Z|X, \boldsymbol{\theta}^{t-1})$$

(2) According above, we can get

$$\begin{aligned}
F(\tilde{P}^t, \boldsymbol{\theta}) &= \mathbb{E}_{\tilde{P}^t} [\log P(X, Z|\boldsymbol{\theta})] + H(\tilde{P}^t) \\
&= \sum_Z P(Z|X, \boldsymbol{\theta}^{t-1}) \log P(X, Z|\boldsymbol{\theta}) + H(\tilde{P}^t) \\
&= Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1}) + H(\tilde{P}^t)
\end{aligned}$$

Then

$$\boldsymbol{\theta}^t = \arg \max_{\boldsymbol{\theta}} F(\tilde{P}^t, \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1})$$

#### 11.4.10.2 The Generalized EM Algorithm(GEM)

In the formulation of the EM algorithm described above,  $\boldsymbol{\theta}^t$  was chosen as the value of  $\boldsymbol{\theta}$  for which  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1})$  was maximized. While this ensures the greatest increase in  $\ell(\boldsymbol{\theta})$ , it is however possible to relax the requirement of maximization to one of simply increasing  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1})$  so that  $Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^{t-1}) \geq Q(\boldsymbol{\theta}^{t-1}, \boldsymbol{\theta}^{t-1})$ . This approach, to simply increase and not necessarily maximize  $Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^{t-1})$  is known as the Generalized Expectation Maximization (GEM) algorithm and is often useful in cases where the maximization is difficult. The convergence of the GEM algorithm is similar to the EM algorithm.

#### 11.4.11 Online EM

#### 11.4.12 Other EM variants \*

### 11.5 Model selection for latent variable models

When using LVMs, we must specify the number of latent variables, which controls the model complexity. In particular, in the case of mixture models, we must specify  $K$ , the number of clusters. Choosing these parameters is an example of model selection. We discuss some approaches below.

### 11.5.1 Model selection for probabilistic models

The optimal Bayesian approach, discussed in Section 5.3, is to pick the model with the largest marginal likelihood,  $K^* = \arg \max_k p(\mathcal{D}|k)$ .

There are two problems with this. First, evaluating the marginal likelihood for LVMs is quite difficult. In practice, simple approximations, such as BIC, can be used (see e.g., (Fraley and Raftery 2002)). Alternatively, we can use the cross-validated likelihood as a performance measure, although this can be slow, since it requires fitting each model  $F$  times, where  $F$  is the number of CV folds.

The second issue is the need to search over a potentially large number of models. The usual approach is to perform exhaustive search over all candidate values of  $K$ . However, sometimes we can set the model to its maximal size, and then rely on the power of the Bayesian Occams razor to kill off unwanted components. An example of this will be shown in Section 21.6.1.6, when we discuss variational Bayes.

An alternative approach is to perform stochastic sampling in the space of models. Traditional approaches, such as (Green 1998, 2003; Lunn et al. 2009), are based on reversible jump MCMC, and use birth moves to propose new centers, and death moves to kill off old centers. However, this can be slow and difficult to implement. A simpler approach is to use a Dirichlet process mixture model, which can be fit using Gibbs sampling, but still allows for an unbounded number of mixture components; see Section 25.2 for details.

Perhaps surprisingly, these sampling-based methods can be faster than the simple approach of evaluating the quality of each  $K$  separately. The reason is that fitting the model for each  $K$  is often slow. By contrast, the sampling methods can often quickly determine that a certain value of  $K$  is poor, and thus they need not waste time in that part of the posterior.

### 11.5.2 Model selection for non-probabilistic methods

What if we are not using a probabilistic model? For example, how do we choose  $K$  for the K-means algorithm? Since this does not correspond to a probability model, there is no likelihood, so none of the methods described above can be used.

An obvious proxy for the likelihood is the **reconstruction error**. Define the squared reconstruction error of a data set  $\mathcal{D}$ , using model complexity  $K$ , as follows:

$$E(\mathcal{D}, K) \triangleq \frac{1}{|\mathcal{D}|} \sum_{i=1}^N \|x_i - \hat{x}_i\|^2 \quad (11.34)$$

In the case of K-means, the reconstruction is given by  $\hat{x}_i = \mu_{z_i}$ , where  $z_i = \arg \min_k \|x_i - \mu_k\|^2$ , as explained in Section 11.4.2.6 TODO.

In supervised learning, we can always use cross validation to select between non-probabilistic models of different complexity, but this is not the case with unsupervised learning. The most common approach is to plot the reconstruction error on the training set versus  $K$ , and to try to identify a **knee** or **kink** in the curve.

## 11.6 Fitting models with missing data

Suppose we want to fit a joint density model by maximum likelihood, but we have holes in our data matrix, due to missing data (usually represented by NaNs). More formally, let  $O_{ij} = 1$  if component  $j$  of data case  $i$  is observed, and let  $O_{ij} = 0$  otherwise. Let  $\mathbf{X}_v = \{x_{ij} : O_{ij} = 1\}$  be the visible data, and  $\mathbf{X}_h = \{x_{ij} : O_{ij} = 0\}$  be the missing or hidden data. Our goal is to compute

$$\hat{\theta} = \arg \max_{\theta} p(\mathbf{X}_v | \theta, \mathbf{O}) \quad (11.35)$$

Under the missing at random assumption (see Section 8.6.2), we have  
TODO

**11.6.1 EM for the MLE of an MVN with missing data**

TODO



## **Chapter 12**

### **Latent linear models**



## **Chapter 13**

### **Sparse linear models**



## Chapter 14

### Kernels

#### 14.1 Introduction

So far in this book, we have been assuming that each object that we wish to classify or cluster or process in anyway can be represented as a fixed-size feature vector, typically of the form  $\mathbf{x}_i \in \mathbb{R}^D$ . However, for certain kinds of objects, it is not clear how to best represent them as fixed-sized feature vectors. For example, how do we represent a text document or protein sequence, which can be of variable length? or a molecular structure, which has complex 3d geometry? or an evolutionary tree, which has variable size and shape?

One approach to such problems is to define a generative model for the data, and use the inferred latent representation and/or the parameters of the model as features, and then to plug these features in to standard methods. For example, in Chapter 28 TODO, we discuss deep learning, which is essentially an unsupervised way to learn good feature representations.

Another approach is to assume that we have some way of measuring the similarity between objects, that doesn't require preprocessing them into feature vector format. For example, when comparing strings, we can compute the edit distance between them. Let  $\kappa(\mathbf{x}, \mathbf{x}') \geq 0$  be some measure of similarity between objects  $\kappa(\mathbf{x}, \mathbf{x}') \in \mathcal{X}$ , we will call  $\kappa$  a **kernel function**. Note that the word kernel has several meanings; we will discuss a different interpretation in Section 14.7.1 TODO.

In this chapter, we will discuss several kinds of kernel functions. We then describe some algorithms that can be written purely in terms of kernel function computations. Such methods can be used when we don't have access to (or choose not to look at) the inside of the objects  $\mathbf{x}$  that we are processing.

#### 14.2 Kernel functions

**Definition 14.1.** A **kernel function**<sup>20</sup> is a real-valued function of two arguments,  $\kappa(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$ . Typically the function is symmetric (i.e.,  $\kappa(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x}', \mathbf{x})$ ), and non-negative (i.e.,  $\kappa(\mathbf{x}, \mathbf{x}') \geq 0$ ).

We give several examples below.

##### 14.2.1 RBF kernels

The **Gaussian kernel** or **squared exponential kernel** (SE kernel) is defined by

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp \left( -\frac{1}{2} (\mathbf{x} - \mathbf{x}')^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \mathbf{x}') \right) \quad (14.1)$$

If  $\boldsymbol{\Sigma}$  is diagonal, this can be written as

---

<sup>20</sup> [http://en.wikipedia.org/wiki/Kernel\\_function](http://en.wikipedia.org/wiki/Kernel_function)

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp \left( -\frac{1}{2} \sum_{j=1}^D \frac{1}{\sigma_j^2} (x_j - x'_j)^2 \right) \quad (14.2)$$

We can interpret the  $\sigma_j$  as defining the **characteristic length scale** of dimension  $j$ . If  $\sigma_j = \infty$ , the corresponding dimension is ignored; hence this is known as the **ARD kernel**. If  $\Sigma$  is spherical, we get the isotropic kernel

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp \left( -\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2} \right) \quad (14.3)$$

Here  $\sigma^2$  is known as the **bandwidth**. Equation (14.4) is an example of a **radial basis function** or **RBF** kernel, since it is only a function of  $\|\mathbf{x} - \mathbf{x}'\|^2$ .

### 14.2.2 TF-IDF kernels

$$\kappa(\mathbf{x}, \mathbf{x}') = \frac{\phi(\mathbf{x})^T \phi(\mathbf{x}')}{\|\phi(\text{vec} \mathbf{x})\|_2 \|\phi(\mathbf{x}')\|_2}, \text{ where } \phi(\mathbf{x}) = \text{tf} - \text{idf}(\mathbf{x}) \quad (14.4)$$

### 14.2.3 Mercer (positive definite) kernels

If the kernel function satisfies the requirement that the **Gram matrix**, defined by

$$\mathbf{K} \triangleq \begin{pmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & & \vdots \\ \kappa(\mathbf{x}_N, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix} \quad (14.5)$$

be positive definite for any set of inputs  $\{\mathbf{x}_i\}_{i=1}^N$ . We call such a kernel a **Mercer kernel**, or **positive definite kernel**.

If the Gram matrix is positive definite, we can compute an eigenvector decomposition of it as follows

$$\mathbf{K} = \mathbf{U}^T \mathbf{\Lambda} \mathbf{U} \quad (14.6)$$

where  $\mathbf{\Lambda}$  is a diagonal matrix of eigenvalues  $\lambda_i > 0$ . Now consider an element of  $\mathbf{K}$ :

$$k_{ij} = (\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}_{:,i})^T (\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}_{:,j}) \quad (14.7)$$

Let us define  $\phi(\mathbf{x}_i) = \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}_{:,i}$ , then we can write

$$k_{ij} = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \quad (14.8)$$

Thus we see that the entries in the kernel matrix can be computed by performing an inner product of some feature vectors that are implicitly defined by the eigenvectors  $\mathbf{U}$ . In general, if the kernel is Mercer, then there exists a function  $\phi$  mapping  $\mathbf{x} \in \mathcal{X}$  to  $\mathbb{R}^D$  such that

$$\kappa(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') \quad (14.9)$$

where  $\phi$  depends on the eigen functions of  $\kappa$  (so  $D$  is a potentially infinite dimensional space).

For example, consider the (non-stationary) **polynomial kernel**  $\kappa(\mathbf{x}, \mathbf{x}') = (\gamma \mathbf{x} \mathbf{x}' + r)^M$ , where  $r > 0$ . One can show that the corresponding feature vector  $\phi(\mathbf{x})$  will contain all terms up to degree  $M$ . For example, if  $M = 2, \gamma = r = 1$  and  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^2$ , we have

$$\begin{aligned}
(\mathbf{x}\mathbf{x}' + 1)^2 &= (1 + x_1x'_1 + x_2 + x'_2)^2 \\
&= 1 + 2x_1x'_1 + 2x_2x'_2 + (x_1x'_1)^2 + (x_2x'_2)^2 + 2x_1x'_1x_2x'_2 \\
&= \phi(\mathbf{x})^T \phi(\mathbf{x}') \quad , \text{ where } \phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)
\end{aligned}$$

In the case of a Gaussian kernel, the feature map lives in an infinite dimensional space. In such a case, it is clearly infeasible to explicitly represent the feature vectors.

In general, establishing that a kernel is a Mercer kernel is difficult, and requires techniques from functional analysis. However, one can show that it is possible to build up new Mercer kernels from simpler ones using a set of standard rules. For example, if  $\kappa_1$  and  $\kappa_2$  are both Mercer, so is  $\kappa(\mathbf{x}, \mathbf{x}') = \kappa_1(\mathbf{x}, \mathbf{x}') + \kappa_2(\mathbf{x}, \mathbf{x}')$ . See e.g., (Schoelkopf and Smola 2002) for details.

### 14.2.4 Linear kernels

$$\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}' \quad (14.10)$$

### 14.2.5 Matern kernels

The **Matern kernel**, which is commonly used in Gaussian process regression (see Section 15.2), has the following form

$$\kappa(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \frac{\sqrt{2\nu}r}{\ell} \right)^\nu K_\nu \frac{\sqrt{2\nu}r}{\ell} \quad (14.11)$$

where  $r = \|\mathbf{x} - \mathbf{x}'\|$ ,  $\nu > 0$ ,  $\ell > 0$ , and  $K_\nu$  is a modified Bessel function. As  $\nu \rightarrow \infty$ , this approaches the SE kernel. If  $\nu = \frac{1}{2}$ , the kernel simplifies to

$$\kappa(r) = \exp(-r/\ell) \quad (14.12)$$

If  $D = 1$ , and we use this kernel to define a Gaussian process (see Chapter 15 TODO), we get the **Ornstein-Uhlenbeck process**, which describes the velocity of a particle undergoing Brownian motion (the corresponding function is continuous but not differentiable, and hence is very jagged).

### 14.2.6 String kernels

Now let  $\phi_s(\mathbf{x})$  denote the number of times that substrings appears in string  $\mathbf{x}$ . We define the kernel between two strings  $\mathbf{x}$  and  $\mathbf{x}'$  as

$$\kappa(\mathbf{x}, \mathbf{x}') = \sum_{s \in \mathcal{A}^*} w_s \phi_s(\mathbf{x}) \phi_s(\mathbf{x}') \quad (14.13)$$

where  $w_s \geq 0$  and  $\mathcal{A}^*$  is the set of all strings (of any length) from the alphabet  $\mathcal{A}$  (this is known as the Kleene star operator). This is a Mercer kernel, and be computed in  $O(|\mathbf{x}| + |\mathbf{x}'|)$  time (for certain settings of the weights  $\{w_s\}$ ) using suffix trees (Leslie et al. 2003; Vishwanathan and Smola 2003; Shawe-Taylor and Cristianini 2004).

There are various cases of interest. If we set  $w_s = 0$  for  $|s| > 1$  we get a bag-of-characters kernel. This defines  $\phi(\mathbf{x})$  to be the number of times each character in  $\mathcal{A}$  occurs in  $\mathbf{x}$ . If we require  $s$  to be bordered by white-space, we get a bag-of-words kernel, where  $\phi(\mathbf{x})$  counts how many times each possible word occurs. Note that this is a very sparse vector, since most words will not be present. If we only consider strings of a fixed length  $k$ , we get the **k-spectrum** kernel. This has been used to classify proteins into SCOP superfamilies (Leslie et al. 2003).

### 14.2.7 Pyramid match kernels

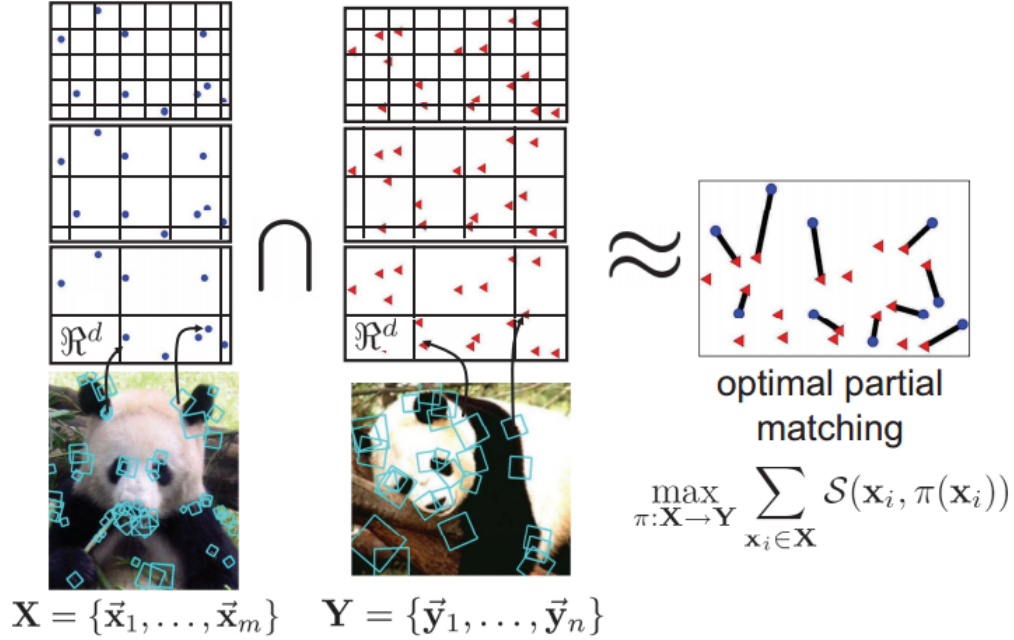


Fig. 14.1: Illustration of a pyramid match kernel computed from two images. Used with kind permission of Kristen Grauman.

### 14.2.8 Kernels derived from probabilistic generative models

Suppose we have a probabilistic generative model of feature vectors,  $p(\mathbf{x}|\boldsymbol{\theta})$ . Then there are several ways we can use this model to define kernel functions, and thereby make the model suitable for discriminative tasks. We sketch two approaches below.

#### 14.2.8.1 Probability product kernels

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \int p(\mathbf{x}|\mathbf{x}_i)^\rho p(\mathbf{x}|\mathbf{x}_j)^\rho d\mathbf{x} \quad (14.14)$$

where  $\rho > 0$ , and  $p(\mathbf{x}|\mathbf{x}_i)$  is often approximated by  $p(\mathbf{x}|\hat{\boldsymbol{\theta}}(\mathbf{x}_i))$ , where  $\hat{\boldsymbol{\theta}}(\mathbf{x}_i)$  is a parameter estimate computed using a single data vector. This is called a **probability product kernel** (Jebara et al. 2004).

Although it seems strange to fit a model to a single data point, it is important to bear in mind that the fitted model is only being used to see how similar two objects are. In particular, if we fit the model to  $\mathbf{x}_i$  and then the model thinks  $\mathbf{x}_j$  is likely, this means that  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are similar. For example, suppose  $p(\mathbf{x}|\boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ , where  $\sigma^2$  is fixed. If  $\rho = 1$ , and we use  $\hat{\boldsymbol{\mu}}(\mathbf{x}_i) = \mathbf{x}_i$  and  $\hat{\boldsymbol{\mu}}(\mathbf{x}_j) = \mathbf{x}_j$ , we find (Jebara et al. 2004, p825) that

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{(4\pi\sigma^2)^{D/2}} \exp\left(-\frac{1}{4\sigma^2} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right) \quad (14.15)$$

which is (up to a constant factor) the RBF kernel.



It turns out that one can compute Equation (14.14) for a variety of generative models, including ones with latent variables, such as HMMs. This provides one way to define kernels on variable length sequences. Furthermore, this technique works even if the sequences are of real-valued vectors, unlike the string kernel in Section 14.2.6. See (Jebara et al. 2004) for further details

### 14.2.8.2 Fisher kernels

A more efficient way to use generative models to define kernels is to use a **Fisher kernel** (Jaakkola and Haussler 1998) which is defined as follows:

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = g(\mathbf{x}_i)^T \mathbf{F}^{-1} g(\mathbf{x}_j) \quad (14.16)$$

where  $g$  is the gradient of the log likelihood, or **score vector**, evaluated at the MLE  $\hat{\theta}$

$$g(\mathbf{x}) \triangleq \frac{d}{d\theta} \log p(\mathbf{x}|\theta)|_{\hat{\theta}} \quad (14.17)$$

and  $\mathbf{F}$  is the Fisher information matrix, which is essentially the Hessian:

$$\mathbf{F} \triangleq \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p(\mathbf{x}|\theta) \right]_{\hat{\theta}} \quad (14.18)$$

Note that  $\hat{\theta}$  is a function of all the data, so the similarity of  $\mathbf{x}_i$  and  $\mathbf{x}_j$  is computed in the context of all the data as well. Also, note that we only have to fit one model.

## 14.3 Using kernels inside GLMs

### 14.3.1 Kernel machines

We define a **kernel machine** to be a GLM where the input feature vector has the form

$$\phi(\mathbf{x}) = (\kappa(\mathbf{x}, \mu_1), \dots, \kappa(\mathbf{x}, \mu_K)) \quad (14.19)$$

where  $\mu_k \in \mathcal{X}$  are a set of  $K$  centroids. If  $\kappa$  is an RBF kernel, this is called an **RBF network**. We discuss ways to choose the  $\mu_k$  parameters below. We will call Equation (14.19) a **kernelised feature vector**. Note that in this approach, the kernel need not be a Mercer kernel.

We can use the kernelized feature vector for logistic regression by defining  $p(y|\mathbf{x}, \theta) = \text{Ber}(y|\mathbf{w}^T \phi(\mathbf{x}))$ . This provides a simple way to define a non-linear decision boundary. For example, see Figure 14.2.

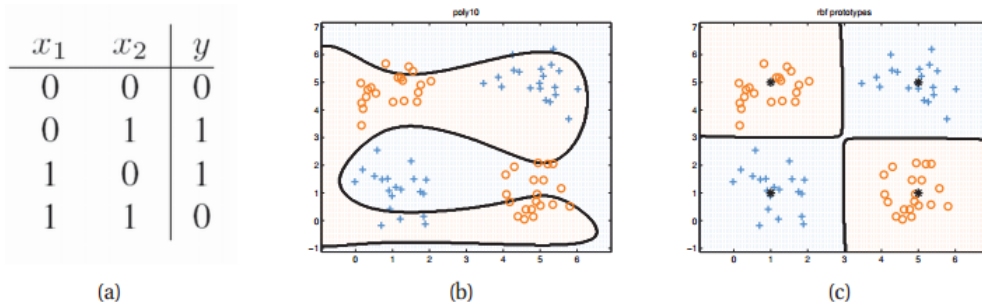


Fig. 14.2: (a) xor truth table. (b) Fitting a linear logistic regression classifier using degree 10 polynomial expansion. (c) Same model, but using an RBF kernel with centroids specified by the 4 black crosses.

We can also use the kernelized feature vector inside a linear regression model by defining  $p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \sigma^2)$ .

### 14.3.2 L1VMs, RVMs, and other sparse vector machines

The main issue with kernel machines is: how do we choose the centroids  $\boldsymbol{\mu}_k$ ? We can use **sparse vector machine**, **L1VM**, **L2VM**, **RVM**, **SVM**.

In Figure 14.3, we compare L2VM, L1VM, RVM and an SVM using the same RBF kernel on a binary classification problem in 2d.

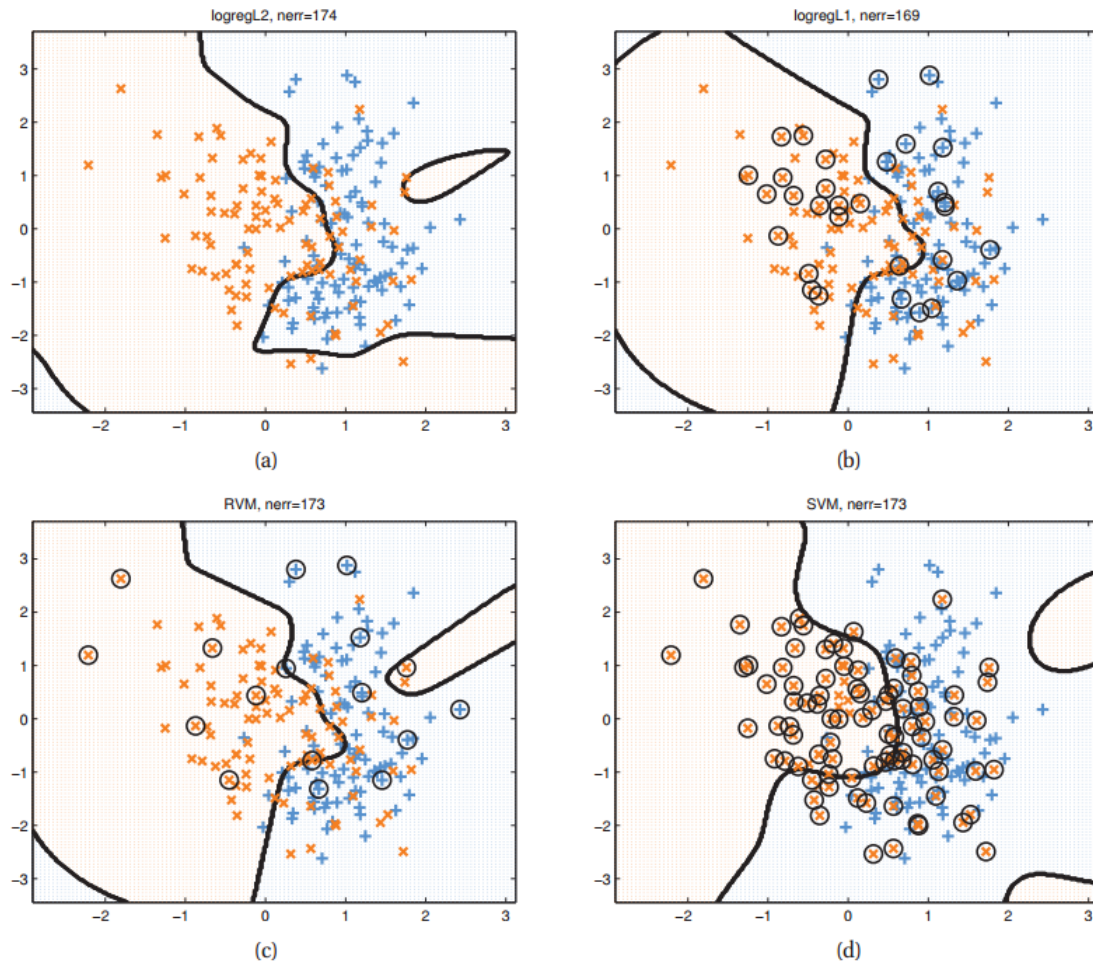


Fig. 14.3: Example of non-linear binary classification using an RBF kernel with bandwidth  $\sigma = 0.3$ . (a) L2VM with  $\lambda = 5$ . (b) L1VM with  $\lambda = 1$ . (c) RVM. (d) SVM with  $C = 1/\lambda$  chosen by cross validation. Black circles denote the support vectors.

In Figure 14.4, we compare L2VM, L1VM, RVM and an SVM using an RBF kernel on a 1d regression problem.

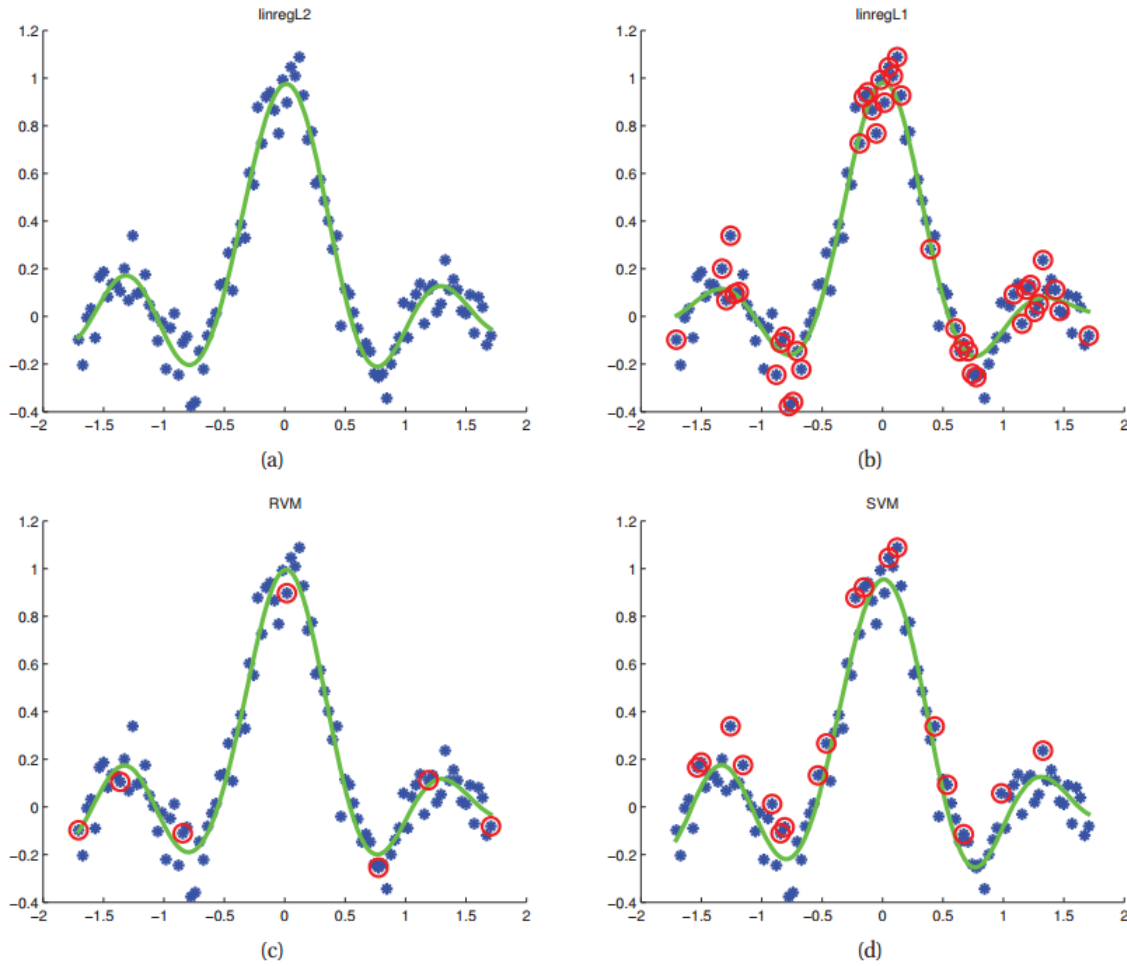


Fig. 14.4: Example of kernel based regression on the noisy sinc function using an RBF kernel with bandwidth  $\sigma = 0.3$ . (a) L2VM with  $\lambda = 0.5$ . (b) L1VM with  $\lambda = 0.5$ . (c) RVM. (d) SVM regression with  $C = 1/\lambda$  chosen by cross validation, and  $\varepsilon = 0.1$  (the default for SVMlight). Red circles denote the retained training exemplars.

## 14.4 The kernel trick

Rather than defining our feature vector in terms of kernels,  $\phi(x) = (\kappa(x, x_1), \dots, \kappa(x, x_N))$ , we can instead work with the original feature vectors  $x$ , but modify the algorithm so that it replaces all inner products of the form  $\langle x_i, x_j \rangle$  with a call to the kernel function,  $\kappa(x_i, x_j)$ . This is called the **kernel trick**. It turns out that many algorithms can be kernelized in this way. We give some examples below. Note that we require that the kernel be a Mercer kernel for this trick to work.

### 14.4.1 Kernelized KNN

The Euclidean distance can be unrolled as

$$\|x_i - x_j\|^2 = \langle x_i, x_i \rangle + \langle x_j, x_j \rangle - 2\langle x_i, x_j \rangle \quad (14.20)$$

then by replacing all  $\langle x_i, x_j \rangle$  with  $\kappa(x_i, x_j)$  we get Kernelized KNN.

## 14.4.2 Kernelized K-medoids clustering

**K-medoids algorithm** is similar to K-means (see Section 11.4.3.6), but instead of representing each clusters centroid by the mean of all data vectors assigned to this cluster, we make each centroid be one of the data vectors themselves. Thus we always deal with integer indexes, rather than data objects.

This algorithm can be kernelized by using Equation (14.20) to replace the distance computation.

## 14.4.3 Kernelized ridge regression

Applying the kernel trick to distance-based methods was straightforward. It is not so obvious how to apply it to parametric models such as ridge regression. However, it can be done, as we now explain. This will serve as a good warm up for studying SVMs.

### 14.4.3.1 The primal problem

we rewrite Equation (7.22) as the following

$$J(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \|\mathbf{w}\|^2 \quad (14.21)$$

and its solution is given by Equation (7.23).

### 14.4.3.2 The dual problem

Equation (14.21) is not yet in the form of inner products. However, using the matrix inversion lemma (Equation 4.107 TODO) we rewrite the ridge estimate as follows

$$\mathbf{w} = \mathbf{X}^T(\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}_N)^{-1}\mathbf{y} \quad (14.22)$$

which takes  $O(N^3 + N^2D)$  time to compute. This can be advantageous if  $D$  is large. Furthermore, we see that we can partially kernelize this, by replacing  $\mathbf{X}\mathbf{X}^T$  with the Gram matrix  $\mathbf{K}$ . But what about the leading  $\mathbf{X}^T$  term?

Let us define the following **dual variables**:

$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1}\mathbf{y} \quad (14.23)$$

Then we can rewrite the **primal variables** as follows

$$\mathbf{w} = \mathbf{X}^T \boldsymbol{\alpha} = \sum_{i=1}^N \alpha_i \mathbf{x}_i \quad (14.24)$$

This tells us that the solution vector is just a linear sum of the  $N$  training vectors. When we plug this in at test time to compute the predictive mean, we get

$$y = f(\mathbf{x}) = \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{x} = \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}) \quad (14.25)$$

So we have successfully kernelized ridge regression by changing from primal to dual variables. This technique can be applied to many other linear models, such as logistic regression.

### 14.4.3.3 Computational cost

The cost of computing the dual variables  $\alpha$  is  $O(N^3)$ , whereas the cost of computing the primal variables  $w$  is  $O(D^3)$ . Hence the kernel method can be useful in high dimensional settings, even if we only use a linear kernel (c.f., the SVD trick in Equation (7.24)). However, prediction using the dual variables takes  $O(ND)$  time, while prediction using the primal variables only takes  $O(D)$  time. We can speedup prediction by making  $\alpha$  sparse, as we discuss in Section 14.5.

## 14.4.4 Kernel PCA

TODO

## 14.5 Support vector machines (SVMs)

In Section ??, we saw one way to derive a sparse kernel machine, namely by using a GLM with kernel basis functions, plus a sparsity-promoting prior such as  $\ell_1$  or ARD. An alternative approach is to change the objective function from negative log likelihood to some other loss function, as we discussed in Section 6.4.5. In particular, consider the  $\ell_2$  regularized empirical risk function

$$J(w, \lambda) = \sum_{i=1}^N L(y_i, \hat{y}_i) + \lambda \|w\|^2, \text{ where } \hat{y}_i = w^T x_i + w_0 \quad (14.26)$$

If  $L$  is quadratic loss, this is equivalent to ridge regression, and if  $L$  is the log-loss defined in Equation (6.3), this is equivalent to logistic regression.

However, if we replace the loss function with some other loss function, to be explained below, we can ensure that the solution is sparse, so that predictions only depend on a subset of the training data, known as **support vectors**. This combination of the kernel trick plus a modified loss function is known as a **support vector machine** or **SVM**.

Note that SVMs are very unnatural from a probabilistic point of view.

- First, they encode sparsity in the loss function rather than the prior.
- Second, they encode kernels by using an algorithmic trick, rather than being an explicit part of the model.
- Finally, SVMs do not result in probabilistic outputs, which causes various difficulties, especially in the multi-class classification setting (see Section 14.5.2.4 TODO for details).

It is possible to obtain sparse, probabilistic, multi-class kernel-based classifiers, which work as well or better than SVMs, using techniques such as the L1VM or RVM, discussed in Section 14.3.2. However, we include a discussion of SVMs, despite their non-probabilistic nature, for two main reasons.

- First, they are very popular and widely used, so all students of machine learning should know about them.
- Second, they have some computational advantages over probabilistic methods in the structured output case; see Section 19.7 TODO.

### 14.5.1 SVMs for classification

#### 14.5.1.1 Primal form

**Representation**

$$\mathcal{H} : y = f(x) = \text{sign}(wx + b) \quad (14.27)$$

**Evaluation**

$$\min_{\mathbf{w}, b} \quad \frac{1}{2} \|\mathbf{w}\|^2 \quad (14.28)$$

$$s.t. \quad y_i(\mathbf{w}\mathbf{x}_i + b) \geq 1, i = 1, 2, \dots, N \quad (14.29)$$

#### 14.5.1.2 Dual form

##### Representation

$$\mathcal{H} : y = f(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^N \alpha_i y_i (\mathbf{x} \cdot \mathbf{x}_i) + b \right) \quad (14.30)$$

##### Evaluation

$$\min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^N \alpha_i \quad (14.31)$$

$$s.t. \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad (14.32)$$

$$\alpha_i \geq 0, i = 1, 2, \dots, N \quad (14.33)$$

#### 14.5.1.3 Primal form with slack variables

##### Representation

$$\mathcal{H} : y = f(\mathbf{x}) = \text{sign}(\mathbf{w}\mathbf{x} + b) \quad (14.34)$$

##### Evaluation

$$\min_{\mathbf{w}, b} \quad C \sum_{i=1}^N \xi_i + \frac{1}{2} \|\mathbf{w}\|^2 \quad (14.35)$$

$$s.t. \quad y_i(\mathbf{w}\mathbf{x}_i + b) \geq 1 - \xi_i, i = 1, 2, \dots, N \quad (14.36)$$

$$\xi_i \geq 0, i = 1, 2, \dots, N \quad (14.37)$$

#### 14.5.1.4 Dual form with slack variables

##### Representation

$$\mathcal{H} : y = f(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^N \alpha_i y_i (\mathbf{x} \cdot \mathbf{x}_i) + b \right) \quad (14.38)$$

##### Evaluation

$$\min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^N \alpha_i \quad (14.39)$$

$$s.t. \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad (14.40)$$

$$0 \leq \alpha_i \leq C, i = 1, 2, \dots, N \quad (14.41)$$

$$\alpha_i = 0 \Rightarrow y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \quad (14.42)$$

$$\alpha_i = C \Rightarrow y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \leq 1 \quad (14.43)$$

$$0 < \alpha_i < C \Rightarrow y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1 \quad (14.44)$$

### 14.5.1.5 Hinge Loss

Linear support vector machines can also be interpreted as hinge loss minimization:

$$\min_{\mathbf{w}, b} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i)) + \lambda \|\mathbf{w}\|^2 \quad (14.45)$$

where  $L(y, f(\mathbf{x}))$  is a **hinge loss function**:

$$L(y, f(\mathbf{x})) = \begin{cases} 1 - yf(\mathbf{x}), & 1 - yf(\mathbf{x}) > 0 \\ 0, & 1 - yf(\mathbf{x}) \leq 0 \end{cases} \quad (14.46)$$

*Proof.* We can write equation (14.45) as equations (14.35)  $\sim$  (14.37).

Define slack variables

$$\xi_i \triangleq 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b), \xi_i \geq 0 \quad (14.47)$$

Then  $\mathbf{w}, b, \xi_i$  satisfy the constraints (14.35) and (14.36). And objective function (14.37) can be written as

$$\min_{\mathbf{w}, b} \sum_{i=1}^N \xi_i + \lambda \|\mathbf{w}\|^2$$

If  $\lambda = \frac{1}{2C}$ , then

$$\min_{\mathbf{w}, b} \frac{1}{C} \left( C \sum_{i=1}^N \xi_i + \frac{1}{2} \|\mathbf{w}\|^2 \right) \quad (14.48)$$

It is equivalent to equation (14.35).

### 14.5.1.6 Optimization

QP, SMO

## 14.5.2 SVMs for regression

### 14.5.2.1 Representation

$$\mathcal{H} : y = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \quad (14.49)$$

### 14.5.2.2 Evaluation

$$J(\mathbf{w}) = C \sum_{i=1}^N L(y_i, f(\mathbf{x}_i)) + \frac{1}{2} \|\mathbf{w}\|^2 \quad (14.50)$$

where  $L(y, f(\mathbf{x}))$  is a **epsilon insensitive loss function**:

$$L(y, f(\mathbf{x})) = \begin{cases} 0 & , |y - f(\mathbf{x})| < \varepsilon \\ |y - f(\mathbf{x})| - \varepsilon & , \text{otherwise} \end{cases} \quad (14.51)$$

and  $C = 1/\lambda$  is a regularization constant.

This objective is convex and unconstrained, but not differentiable, because of the absolute value function in the loss term. As in Section 13.4 TODO, where we discussed the lasso problem, there are several possible algorithms we

could use. One popular approach is to formulate the problem as a constrained optimization problem. In particular, we introduce **slack variables** to represent the degree to which each point lies outside the tube:

$$\begin{aligned} y_i &\leq f(\mathbf{x}_i) + \varepsilon + \xi_i^+ \\ y_i &\geq f(\mathbf{x}_i) - \varepsilon - \xi_i^- \end{aligned}$$

Given this, we can rewrite the objective as follows:

$$J(\mathbf{w}) = C \sum_{i=1}^N (\xi_i^+ + \xi_i^-) + \frac{1}{2} \|\mathbf{w}\|^2 \quad (14.52)$$

This is a standard quadratic problem in  $2N + D + 1$  variables.

### 14.5.3 Choosing $C$

SVMs for both classification and regression require that you specify the kernel function and the parameter  $C$ . Typically  $C$  is chosen by cross-validation. Note, however, that  $C$  interacts quite strongly with the kernel parameters. For example, suppose we are using an RBF kernel with precision  $\gamma = \frac{1}{2\sigma^2}$ . If  $\gamma = 5$ , corresponding to narrow kernels, we need heavy regularization, and hence small  $C$  (so  $\lambda = 1/C$  is big). If  $\gamma = 1$ , a larger value of  $C$  should be used. So we see that  $\gamma$  and  $C$  are tightly coupled. This is illustrated in Figure 14.5, which shows the CV estimate of the 0-1 risk as a function of  $C$  and  $\gamma$ .

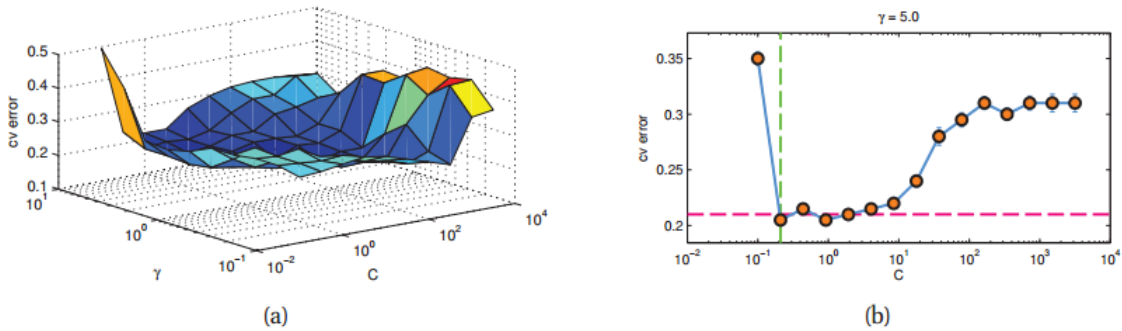


Fig. 14.5: (a) A cross validation estimate of the 0-1 error for an SVM classifier with RBF kernel with different precisions  $\gamma = 1/(2\sigma^2)$  and different regularizer  $\gamma = 1/C$ , applied to a synthetic data set drawn from a mixture of 2 Gaussians. (b) A slice through this surface for  $\gamma = 5$ . The red dotted line is the Bayes optimal error, computed using Bayes rule applied to the model used to generate the data. Based on Figure 12.6 of (Hastie et al. 2009).

The authors of libsvm recommend (Hsu et al. 2009) using CV over a 2d grid with values  $C \in \{2^5, 2^3, \dots, 2^{15}\}$  and  $\gamma \in \{2^{15}, 2^{13}, \dots, 2^3\}$ . In addition, it is important to standardize the data first, for a spherical Gaussian kernel to make sense.

To choose  $C$  efficiently, one can develop a path following algorithm in the spirit of lars (Section 13.3.4 TODO). The basic idea is to start with  $\lambda$  large, so that the margin  $1/\|\mathbf{w}(\lambda)\|$  is wide, and hence all points are inside of it and have  $\alpha_i = 1$ . By slowly decreasing  $\lambda$ , a small set of points will move from inside the margin to outside, and their  $\alpha_i$  values will change from 1 to 0, as they cease to be support vectors. When  $\lambda$  is maximal, the function is completely smoothed, and no support vectors remain. See (Hastie et al. 2004) for the details.



### 14.5.4 A probabilistic interpretation of SVMs

TODO see MLAPP Section 14.5.5

### 14.5.5 Summary of key points

Summarizing the above discussion, we recognize that SVM classifiers involve three key ingredients: the kernel trick, sparsity, and the large margin principle. The kernel trick is necessary to prevent underfitting, i.e., to ensure that the feature vector is sufficiently rich that a linear classifier can separate the data. (Recall from Section 14.2.3 that any Mercer kernel can be viewed as implicitly defining a potentially high dimensional feature vector.) If the original features are already high dimensional (as in many gene expression and text classification problems), it suffices to use a linear kernel,  $\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$ , which is equivalent to working with the original features.

The sparsity and large margin principles are necessary to prevent overfitting, i.e., to ensure that we do not use all the basis functions. These two ideas are closely related to each other, and both arise (in this case) from the use of the hinge loss function. However, there are other methods of achieving sparsity (such as  $\ell_1$ ), and also other methods of maximizing the margin (such as boosting). A deeper discussion of this point takes us outside of the scope of this book. See e.g., (Hastie et al. 2009) for more information.

## 14.6 Comparison of discriminative kernel methods

We have mentioned several different methods for classification and regression based on kernels, which we summarize in Table 14.1. (GP stands for Gaussian process, which we discuss in Chapter 15 TODO.) The columns have the following meaning:

- Optimize  $\mathbf{w}$ : a key question is whether the objective  $J(\mathbf{w} = -\log p(\mathcal{D}|\mathbf{w}) - \log p(\mathbf{w}))$  is convex or not. L2VM, L1VM and SVMs have convex objectives. RVMs do not. GPs are Bayesian methods that do not perform parameter estimation.
- Optimize kernel: all the methods require that one tune the kernel parameters, such as the bandwidth of the RBF kernel, as well as the level of regularization. For methods based on Gaussians, including L2VM, RVMs and GPs, we can use efficient gradient based optimizers to maximize the marginal likelihood. For SVMs, and L1VM, we must use cross validation, which is slower (see Section 14.5.3).
- Sparse: L1VM, RVMs and SVMs are sparse kernel methods, in that they only use a subset of the training examples. GPs and L2VM are not sparse: they use all the training examples. The principle advantage of sparsity is that prediction at test time is usually faster. In addition, one can sometimes get improved accuracy.
- Probabilistic: All the methods except for SVMs produce probabilistic output of the form  $p(y|\mathbf{x})$ . SVMs produce a confidence value that can be converted to a probability, but such probabilities are usually very poorly calibrated (see Section 14.5.2.3 TODO).
- Multiclass: All the methods except for SVMs naturally work in the multiclass setting, by using a multinoulli output instead of Bernoulli. The SVM can be made into a multiclass classifier, but there are various difficulties with this approach, as discussed in Section 14.5.2.4 TODO.
- Mercer kernel: SVMs and GPs require that the kernel is positive definite; the other techniques do not.

Method	Opt. $\mathbf{w}$	Opt. kernel	Sparse	Prob.	Multiclass	Non-Mercer	Section
L2VM	Convex	EB	No	Yes	Yes	Yes	14.3.2
L1VM	Convex	CV	Yes	Yes	Yes	Yes	14.3.2
RVM	Not convex	EB	Yes	Yes	Yes	Yes	14.3.2
SVM	Convex	CV	Yes	No	Indirectly	No	14.5
GP	N/A	EB	No	Yes	Yes	No	15

---

Table 14.1: Comparison of various kernel based classifiers. EB = empirical Bayes, CV = cross validation. See text for details

## **14.7 Kernels for building generative models**

TODO

## Chapter 15

### Perceptron

#### 15.1 Representation

$$\mathcal{H} : y = f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b) \quad (15.1)$$

where  $\text{sign}(x) = \begin{cases} +1, & x \geq 0 \\ -1, & x < 0 \end{cases}$ , see Fig. 15.1<sup>21</sup>.

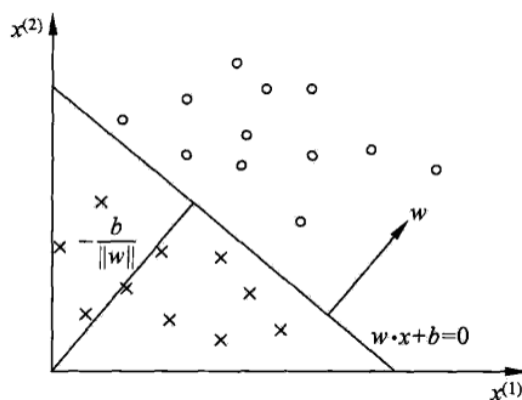


Fig. 15.1: Perceptron

#### 15.2 Evaluation

$$L(\mathbf{w}, b) = -y_i(\mathbf{w}^T \mathbf{x}_i + b) \quad (15.2)$$

$$R_{emp}(f) = -\sum_i y_i(\mathbf{w}^T \mathbf{x}_i + b) \quad (15.3)$$

$$(15.4)$$

<sup>21</sup> <https://en.wikipedia.org/wiki/Perceptron>

## 15.3 Optimization

### 15.3.1 Primal form

Stochastic gradient descent, the pseudo code is as follows:

```

 $w \leftarrow 0; b \leftarrow 0; k \leftarrow 0;$ 
while no mistakes made within the for loop do
  for  $i \leftarrow 1$  to  $N$  do
    if  $y_i(w \cdot x_i + b) \leq 0$  then
       $w \leftarrow w + \eta y_i x_i;$ 
       $b \leftarrow b + \eta y_i;$ 
       $k \leftarrow k + 1;$ 
    end
  end
end

```

**Algorithm 3:** Perceptron learning algorithm, primal form

#### 15.3.1.1 Convergency

**Theorem 15.1.** (Novikoff) If training data set  $\mathcal{D}$  is linearly separable, then

1. There exists a hyperplane denoted as  $\hat{w}_{opt} \cdot x + b_{opt} = 0$  which can correctly separate all samples, and

$$\exists \gamma > 0, \forall i, y_i(w_{opt} \cdot x_i + b_{opt}) \geq \gamma \quad (15.5)$$

- 2.

$$k \leq \left(\frac{R}{\gamma}\right)^2, \text{ where } R = \max_{1 \leq i \leq N} \|\hat{x}_i\| \quad (15.6)$$

*Proof.* (1) let  $\gamma = \min_i y_i(w_{opt} \cdot x_i + b_{opt})$ , then we get  $y_i(w_{opt} \cdot x_i + b_{opt}) \geq \gamma$ .

(2) The algorithm start from  $\hat{x}_0 = 0$ , if a instance is misclassified, then update the weight. Let  $\hat{w}_{k-1}$  denotes the extended weight before the k-th misclassified instance, then we can get

$$y_i(\hat{w}_{k-1} \cdot \hat{x}_i) = y_i(w_{k-1} \cdot x_i + b_{k-1}) \leq 0 \quad (15.7)$$

$$\hat{w}_k = \hat{w}_{k-1} + \eta y_i \hat{x}_i \quad (15.8)$$

We could infer the following two equations, the proof procedure are omitted.

1.  $\hat{w}_k \cdot \hat{w}_{opt} \geq k\eta\gamma$
2.  $\|\hat{w}_k\|^2 \leq k\eta^2 R^2$

From above two equations we get

$$\begin{aligned}
 k\eta\gamma &\leq \hat{w}_k \cdot \hat{w}_{opt} \leq \|\hat{w}_k\| \|\hat{w}_{opt}\| \leq \sqrt{k}\eta R \\
 k^2\gamma^2 &\leq kR^2 \\
 \text{i.e. } k &\leq \left(\frac{R}{\gamma}\right)^2
 \end{aligned}$$

### 15.3.2 Dual form

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \quad (15.9)$$

$$b = \sum_{i=1}^N \alpha_i y_i \quad (15.10)$$

$$f(\mathbf{x}) = \text{sign} \left( \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \cdot \mathbf{x} + b \right) \quad (15.11)$$

```

 $\alpha \leftarrow 0$ ;  $b \leftarrow 0$ ;  $k \leftarrow 0$ ;
while no mistakes made within the for loop do
  for  $i \leftarrow 1$  to  $N$  do
    if  $y_i \left( \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \cdot \mathbf{x}_i + b \right) \leq 0$  then
       $\alpha \leftarrow \alpha + \eta$ ;
       $b \leftarrow b + \eta y_i$ ;
       $k \leftarrow k + 1$ ;
    end
  end
end

```

**Algorithm 4:** Perceptron learning algorithm, dual form



## Chapter 16

### K-Nearest Neighbors

#### 16.1 Representation

$$y = f(\mathbf{x}) = \arg \min_c \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} \mathbb{I}(y_i = c) \quad (16.1)$$

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.

#### 16.2 Evaluation

No training is needed.

#### 16.3 Optimization

No training is needed.





## Chapter 17

### K-Means Clustering

#### 17.1 Representation

$$y_j = k \text{ if } \|\mathbf{x}_j - \boldsymbol{\mu}_k\|_2^2 \text{ is minimal} \quad (17.1)$$

where  $\boldsymbol{\mu}_k$  is the centroid of cluster k.

#### 17.2 Evaluation

$$\arg \min_{\boldsymbol{\mu}} \sum_{j=1}^N \sum_{k=1}^K \gamma_{jk} \|\mathbf{x}_j - \boldsymbol{\mu}_k\|_2^2 \quad (17.2)$$

The hidden variable is  $\gamma_{jk}$ , which's meaning is:

$$\gamma_{jk} = \begin{cases} 1, & \text{if } \|\mathbf{x}_j - \boldsymbol{\mu}_k\|_2 \text{ is minimal for } \boldsymbol{\mu}_k \\ 0, & \text{otherwise} \end{cases}$$

#### 17.3 Optimization

E-Step:

$$\gamma_{jk}^{(i+1)} = \begin{cases} 1, & \text{if } \|\mathbf{x}_j - \boldsymbol{\mu}_k^{(i)}\|_2 \text{ is minimal for } \boldsymbol{\mu}_k^{(i)} \\ 0, & \text{otherwise} \end{cases} \quad (17.3)$$

M-Step:

$$\boldsymbol{\mu}_k^{(i+1)} = \frac{\sum_{j=1}^N \gamma_{jk}^{(i+1)} \mathbf{x}_j}{\sum \gamma_{jk}^{(i+1)}} \quad (17.4)$$

## 17.4 Tricks

### 17.4.1 Choosing $k$

### 17.4.2 Choosing the initial centroids(seeds)

#### 17.4.2.1 K-means++

The intuition that spreading out the  $k$  initial cluster centers is a good thing is behind this approach: the first cluster center is chosen uniformly at random from the data points that are being clustered, after which each subsequent cluster center is chosen from the remaining data points with probability proportional to its squared distance from the point's closest existing cluster center<sup>22</sup>.

The exact algorithm is as follows:

1. Choose one center uniformly at random from among the data points.
2. For each data point  $x$ , compute  $D(x)$ , the distance between  $x$  and the nearest center that has already been chosen.
3. Choose one new data point at random as a new center, using a weighted probability distribution where a point  $x$  is chosen with probability proportional to  $D(x)^2$ .
4. Repeat Steps 2 and 3 until  $k$  centers have been chosen.
5. Now that the initial centers have been chosen, proceed using standard k-means clustering.

## 17.5 Reference

1. cheat-sheet: Algorithm for supervised and unsupervised learning by Emanuel Ferm <http://t.cn/hD0Stf>

---

<sup>22</sup> <http://en.wikipedia.org/wiki/K-means++>

## **Chapter 18**

### **Decision Tree**



## Chapter 19

### AdaBoost

#### 19.1 Representation

$$y = \text{sign}(f(\mathbf{x})) = \text{sign} \left( \sum_{i=1}^m \alpha_i G_i(\mathbf{x}) \right) \quad (19.1)$$

where  $G_i(\mathbf{x})$  are sub classifiers.

#### 19.2 Evaluation

$L(y, f(\mathbf{x})) = \exp[-yf(\mathbf{x})]$  i.e., exponential loss function

$$(\alpha_i, G_i(x)) = \arg \min_{\alpha, G} \sum_{i=1}^N \exp[-y_i(f_{i-1}(\mathbf{x}_i) + \alpha G(\mathbf{x}_i))] \quad (19.2)$$

Define  $\bar{w}_{mi} = \exp[-y_i(f_{i-1}(\mathbf{x}_i))]$ , which is constant w.r.t.  $\alpha, G$

$$(\alpha_i, G_i(x)) = \arg \min_{\alpha, G} \sum_{i=1}^N \bar{w}_{mi} \exp(-y_i \alpha G(\mathbf{x}_i)) \quad (19.3)$$

#### 19.3 Optimization

##### 19.3.1 Input

$\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ , where  $\mathbf{x}_i \in \mathbb{R}^D$ ,  $y_i \in \{-1, +1\}$   
Weak classifiers  $\{G_1, G_2, \dots, G_m\}$

##### 19.3.2 Output

Final classifier:  $G(\mathbf{x})$

### 19.3.3 Algorithm

1. Initialize the weights' distribution of training data(when  $m = 1$ )

$$\mathcal{D}_1 = (w_{11}, w_{12}, \dots, w_{1n}) = (\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N}), i = 1, 2, \dots, N$$

2. Iterate over  $m = 1, 2, \dots, M$

- (a) Use training data with current weights' distribution  $\mathcal{D}_m$  to get a classifier  $G_m(\mathbf{x})$
- (b) Compute the error rate of  $G_m(\mathbf{x})$  over the training data

$$e_m = P(G_m(\mathbf{x}_i) \neq y_i) = \sum_{i=1}^N w_{mi} \mathbb{I}(G_m(\mathbf{x}_i) \neq y_i) \quad (19.4)$$

- (c) Compute the coefficient of classifier  $G_m(\mathbf{x})$

$$\alpha_m = \frac{1}{2} \log \frac{1 - e_m}{e_m} \quad (19.5)$$

- (d) Update the weights' distribution of training data

$$w_{m+1,i} = \frac{w_{mi}}{Z_m} \exp(-\alpha_m y_i G_m(\mathbf{x}_i)) \quad (19.6)$$

where  $Z_m$  is the normalizing constant

$$Z_m = \sum_{i=1}^N w_{mi} \exp(-\alpha_m y_i G_m(\mathbf{x}_i)) \quad (19.7)$$

3. Ensemble  $M$  weak classifiers

$$G(\mathbf{x}) = \text{sign} f(\mathbf{x}) = \text{sign} \left[ \sum_{m=1}^M \alpha_m G_m(\mathbf{x}) \right] \quad (19.8)$$

## 19.4 The upper bound of the training error of AdaBoost

**Theorem 19.1.** *The upper bound of the training error of AdaBoost is*

$$\frac{1}{N} \sum_{i=1}^N \mathbb{I}(G(\mathbf{x}_i) \neq y_i) \leq \frac{1}{N} \sum_{i=1}^N \exp(-y_i f(\mathbf{x}_i)) = \prod_{m=1}^M Z_m \quad (19.9)$$

*Note: the following equation would help proof this theorem*

$$w_{mi} \exp(-\alpha_m y_i G_m(\mathbf{x}_i)) = Z_m w_{m+1,i} \quad (19.10)$$

## **Chapter 20**

### **Hidden markov Model**





## **Chapter 21**

# **Conditional Random Field**



## Appendix A

### Optimization methods

#### A.1 Convexity

**Definition A.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

$$\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \mathcal{S}, \forall \lambda \in [0, 1] \quad (\text{A.1})$$

**Definition A.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

$$f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) \leq \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2) \quad (\text{A.2})$$

**Definition A.3.** A function  $f(\mathbf{x})$  is said to be **strictly convex** if the inequality is strict

$$f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) < \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2) \quad (\text{A.3})$$

**Definition A.4.** A function  $f(\mathbf{x})$  is said to be (strictly) **concave** if  $-f(\mathbf{x})$  is (strictly) convex.

**Theorem A.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 A.1.**  $\log(x)$  is strictly concave 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.

```

input : Training data  $\mathcal{D} = \{(\mathbf{x}_i, y_i) | i = 1 : N\}$ 
output: A linear model:  $y_i = \boldsymbol{\theta}^T \mathbf{x}$ 
 $\mathbf{w} \leftarrow 0$ ;  $b \leftarrow 0$ ;  $k \leftarrow 0$ ;
while no mistakes made within the for loop do
  for  $i \leftarrow 1$  to  $N$  do
    if  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \leq 0$  then
       $\mathbf{w} \leftarrow \mathbf{w} + \eta y_i \mathbf{x}_i$ ;
       $b \leftarrow b + \eta y_i$ ;
       $k \leftarrow k + 1$ ;
    end
  end
end

```

**Algorithm 5:** Stochastic gradient descent

## A.2 Gradient descent

### A.2.1 Stochastic gradient descent

### A.2.2 Batch gradient descent

### A.2.3 Line search

The **line search**<sup>1</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 A.2), Newton's method(Section A.4) and Quasi-Newton method(Section A.5). The step size can be determined either exactly or inexactly.

### A.2.4 Momentum term

## A.3 Lagrange duality

### A.3.1 Primal form

Consider the following, which we'll call the **primal** optimization problem:

$$\text{xyz} \tag{A.4}$$

### A.3.2 Dual form

## A.4 Newton's method

<sup>1</sup> [http://en.wikipedia.org/wiki/Line\\_search](http://en.wikipedia.org/wiki/Line_search)

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \mathbf{g}_k^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^T \mathbf{H}_k (\mathbf{x} - \mathbf{x}_k)$$

$$\text{where } \mathbf{g}_k \triangleq \mathbf{g}(\mathbf{x}_k) = f'(\mathbf{x}_k)$$

$$\mathbf{H}_k \triangleq \mathbf{H}(\mathbf{x}_k), \mathbf{H}(\mathbf{x}) \triangleq \left[ \frac{\partial^2 f}{\partial x_i \partial x_j} \right]_{D \times D} \quad (\text{Hessian matrix})$$

$$f'(\mathbf{x}) = \mathbf{g}_k + \mathbf{H}_k (\mathbf{x} - \mathbf{x}_k) = 0 \Rightarrow \quad (\text{A.5})$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}_k^{-1} \mathbf{g}_k \quad (\text{A.6})$$

Initialize  $\mathbf{x}_0$

**while** (!convergency) **do**

    Evaluate  $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$

    Evaluate  $\mathbf{H}_k = \nabla^2 f(\mathbf{x}_k)$

$\mathbf{d}_k = -\mathbf{H}_k^{-1} \mathbf{g}_k$

    Use line search to find step size  $\eta_k$  along  $\mathbf{d}_k$

$\mathbf{x}_{k+1} = \mathbf{x}_k + \eta_k \mathbf{d}_k$

**end**

**Algorithm 6:** Newtons method for minimizing a strictly convex function

## A.5 Quasi-Newton method

From Equation (A.5) we can infer out the **quasi-Newton condition** as follows:

$$\begin{aligned} f'(\mathbf{x}) - \mathbf{g}_k &= \mathbf{H}_k (\mathbf{x} - \mathbf{x}_k) \\ \mathbf{g}_{k-1} - \mathbf{g}_k &= \mathbf{H}_k (\mathbf{x}_{k-1} - \mathbf{x}_k) \Rightarrow \\ \mathbf{g}_k - \mathbf{g}_{k-1} &= \mathbf{H}_k (\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{g}_{k+1} - \mathbf{g}_k &= \mathbf{H}_{k+1} (\mathbf{x}_{k+1} - \mathbf{x}_k) \quad (\text{quasi-Newton condition}) \end{aligned} \quad (\text{A.7})$$

The idea is to replace  $\mathbf{H}_k^{-1}$  with a approximation  $\mathbf{B}_k$ , which satisfies the following properties:

1.  $\mathbf{B}_k$  must be symmetric
2.  $\mathbf{B}_k$  must satisfies the quasi-Newton condition, i.e.,  $\mathbf{g}_{k+1} - \mathbf{g}_k = \mathbf{B}_{k+1} (\mathbf{x}_{k+1} - \mathbf{x}_k)$ .

Let  $\mathbf{y}_k = \mathbf{g}_{k+1} - \mathbf{g}_k$ ,  $\delta_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ , then

$$\mathbf{B}_{k+1} \mathbf{y}_k = \delta_k \quad (\text{A.8})$$

3. Subject to the above,  $\mathbf{B}_k$  should be as close as possible to  $\mathbf{B}_{k-1}$ .

Note that we did not require that  $\mathbf{B}_k$  be positive definite. That is because we can show that it must be positive definite if  $\mathbf{B}_{k-1}$  is. Therefore, as long as the initial Hessian approximation  $\mathbf{B}_0$  is positive definite, all  $\mathbf{B}_k$  are, by induction.

### A.5.1 DFP

Updating rule:

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \mathbf{P}_k + \mathbf{Q}_k \quad (\text{A.9})$$

From Equation (A.8) we can get

$$\mathbf{B}_{k+1} \mathbf{y}_k = \mathbf{B}_k \mathbf{y}_k + \mathbf{P}_k \mathbf{y}_k + \mathbf{Q}_k \mathbf{y}_k = \delta_k$$

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

$$P_k = \frac{\delta_k \delta_k^T}{\delta_k^T y_k} \quad (\text{A.10})$$

$$Q_k = -\frac{B_k y_k y_k^T B_k}{y_k^T B_k y_k} \quad (\text{A.11})$$

### A.5.2 BFGS

Use  $B_k$  as a approximation to  $H_k$ , then the quasi-Newton condition becomes

$$B_{k+1} \delta_k = y_k \quad (\text{A.12})$$

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

$$P_k = \frac{y_k y_k^T}{y_k^T \delta_k} \quad (\text{A.13})$$

$$Q_k = -\frac{B_k \delta_k \delta_k^T B_k}{\delta_k^T B_k \delta_k} \quad (\text{A.14})$$

### A.5.3 Broyden

Broyden's algorithm is a linear combination of DFP and BFGS.

# Glossary

**feature vector** A feature vector to represent one data.

**loss function** a function that maps an event onto a real number intuitively representing some "cost" associated with the event.

**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.