

# Machine Learning Notebook

Cong Bao

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	About this Notebook . . . . .	1
1.2	Policy of Use . . . . .	1
<b>2</b>	<b>Mathematics Basics</b>	<b>2</b>
2.1	Probability . . . . .	2
2.1.1	Basic Rules . . . . .	2
2.1.2	Common Probability Distributions . . . . .	4
2.2	Linear Algebra . . . . .	6
2.2.1	Vectors . . . . .	6
2.2.2	Matrices . . . . .	8
2.3	Calculus . . . . .	8
2.3.1	Differentiation & Integration . . . . .	8
2.3.2	Multivariate Calculus . . . . .	8
2.4	Informatics . . . . .	8
2.4.1	Entropy . . . . .	8
2.5	Optimization . . . . .	8
2.5.1	One-dimensional Minimization . . . . .	8
2.5.2	Gradient Descent . . . . .	8
2.5.3	Quadratic Functions . . . . .	8
2.5.4	General Functions . . . . .	8
2.5.5	Optimization with Constraints . . . . .	8
<b>3</b>	<b>Machine Learning Basics</b>	<b>9</b>
3.1	Regularization . . . . .	9
3.1.1	Under-fitting & Over-fitting . . . . .	9
3.1.2	Bias & Variance . . . . .	9
3.1.3	Vector Norm . . . . .	9
3.1.4	Penalize Complexity . . . . .	10

3.2	Cross-Validation . . . . .	10
3.3	Bayesian Learning . . . . .	10
3.3.1	Bayes' Rule Terminology . . . . .	10
3.3.2	Maximum Likelihood . . . . .	11
3.3.3	Maximum a Posterior (MAP) . . . . .	11
3.3.4	Bayesian Approach . . . . .	11
3.3.5	Example: Univariate Normal Distribution . . . . .	12
3.3.6	Example: Categorical Distribution . . . . .	14
3.4	Machine Learning Models . . . . .	16
3.4.1	Learning and Inference . . . . .	16
3.4.2	Three Types of Model . . . . .	17
3.4.3	Example: Regression . . . . .	18
3.4.4	Example: Classification . . . . .	19
3.5	Overview of Common Algorithms . . . . .	20
<b>4</b>	<b>Linear Regression</b>	<b>21</b>
4.1	Basic Model . . . . .	21
4.2	Bayesian Regression . . . . .	21
4.3	Non-linear Regression . . . . .	22
4.4	Kernel Trick & Gaussian Processes . . . . .	22
4.5	Sparse Linear Regression . . . . .	22
4.6	Dual Linear Regression . . . . .	22
4.7	Relevance Vector Regression . . . . .	22
<b>5</b>	<b>Logistic Regression</b>	<b>23</b>
5.1	Logistic Regression . . . . .	23
5.2	Non-linear Logistic Regression . . . . .	23
5.3	Kernel Trick & Gaussian Process Classification . . . . .	23
5.4	Multi-class Classification . . . . .	23
<b>6</b>	<b>Support Vector Machines</b>	<b>24</b>
6.1	Geometric Margins . . . . .	24
6.2	Primal & Dual Problems . . . . .	24
6.3	Support Vectors . . . . .	24
6.4	Slack Variables . . . . .	24
6.5	Hinge Loss . . . . .	24
6.6	Non-linear SVMs . . . . .	24
6.7	Kernel Trick . . . . .	24

<b>7</b>	<b>EM Algorithm</b>	<b>25</b>
7.1	Expectation Maximization . . . . .	25
7.2	Example: Mixture of Gaussians . . . . .	25
7.3	Example: t-distributions . . . . .	25
7.4	Example: Factor Analysis . . . . .	25
<b>8</b>	<b>Boosting</b>	<b>26</b>
8.1	Ensemble Methods . . . . .	26
8.2	Bagging . . . . .	26
8.3	Boosting . . . . .	26
8.4	Adaboost . . . . .	26
<b>9</b>	<b>Decision Tree &amp; Random Forest</b>	<b>27</b>
9.1	CART . . . . .	27
9.2	ID3 . . . . .	27
9.3	C4.5 . . . . .	27
9.4	Random Forest . . . . .	27
<b>10</b>	<b>Graphical Models &amp; Markov Network</b>	<b>28</b>
10.1	Graph Definitions . . . . .	28
10.1.1	Graph . . . . .	28
10.1.2	Directed Graph . . . . .	28
10.1.3	Undirected Graph . . . . .	29
10.1.4	Connectivity . . . . .	29
10.1.5	Connectedness . . . . .	29
10.2	Belief Networks . . . . .	30
10.2.1	Definition . . . . .	30
10.2.2	Uncertain Evidence . . . . .	30
10.2.3	Independence . . . . .	31
10.2.4	General Rule for Independence in Belief Networks . . . . .	32
10.2.5	Markov Equivalence . . . . .	33
10.3	Markov Networks . . . . .	33
10.3.1	Definition . . . . .	33
10.3.2	Examples . . . . .	33
10.3.3	Independence . . . . .	34
10.3.4	Expressiveness of Markov and Belief Networks . . . . .	34
10.3.5	Factor Graphs . . . . .	34
10.4	Markov Chains . . . . .	34

10.5 Hidden Markov Models . . . . .	34
<b>A Statistical Assessment</b>	<b>35</b>
A.1 Hypothesis Testing . . . . .	35
A.1.1 Testing Basics . . . . .	35
A.1.2 Testing Procedure . . . . .	35
A.1.3 Power Investigation . . . . .	36
A.1.4 Useful Tests . . . . .	36
A.2 Confidence Intervals . . . . .	36
A.3 Bootstrap . . . . .	36

# Chapter 1

## Introduction

### 1.1 About this Notebook

### 1.2 Policy of Use

# Chapter 2

## Mathematics Basics

### 2.1 Probability

#### 2.1.1 Basic Rules

**Three Axioms of Probability** Let  $\Omega$  be a sample space. A probability assigns a real number  $P(X)$  to each event  $X \subseteq \Omega$  in such a way that

1.  $P(X) \geq 0, \forall X$
2. If  $X_1, X_2, \dots$  are pairwise disjoint events ( $X_1 \cap X_2 = \emptyset, i \neq j, i, j = 1, 2, \dots$ ), then  $P(\bigcup_{i=1}^{\infty} X_i) = \sum_{i=1}^{\infty} P(X_i)$ . (This property is called countable additivity.)
3.  $P(\Omega) = 1$

**Joint Probability** The probability both event A and B occur.  $P(X, Y) = P(X \cap Y)$ .

**Marginalization** The probability distribution of any variable in a joint distribution can be recovered by integrating (or summing) over the other variables.

1. For continuous r.v.  $P(x) = \int P(x, y) dy ; P(y) = \int P(x, y) dx$ .
2. For discrete r.v.  $P(x) = \sum_y P(x, y) ; P(y) = \sum_x P(x, y)$ .
3. For mixed r.v.  $P(x, y) = \sum_w \int P(w, x, y, z) dz$ , where  $w$  is discrete and  $z$  is continuous.

**Conditional Probability**  $P(X = x|Y = y)$  is the probability  $X = x$  occurs given the knowledge  $Y = y$  occurs. Conditional probability can be extracted from joint probability that

$$P(x|y = y^*) = \frac{P(x, y = y^*)}{\int P(x, y = y^*) dx} = \frac{P(x, y = y^*)}{P(y = y^*)}$$

Usually, the formula is written as  $P(x|y) = \frac{P(x, y)}{P(y)}$ .

**Product Rule** The formula can be rearranged as  $P(x, y) = P(x|y) P(y) = P(y|x) P(x)$ .  
In case of multiple variables

$$\begin{aligned} P(w, x, y, z) &= P(w, x, y|z) P(z) \\ &= P(w, x|y, z) P(y|z) P(z) \\ &= P(w|x, y, z) P(x|y, z) P(y|z) P(z) \end{aligned}$$

**Independence** If two variables  $x$  and  $y$  are independent, then r.v.  $x$  tells nothing about r.v.  $y$  (and vice-versa)

$$\begin{aligned} P(x|y) &= P(x) \\ P(y|x) &= P(y) \\ P(x, y) &= P(x) P(y) \end{aligned}$$

**Baye's Rule** By rearranging formula in Product Rule, we have

$$\begin{aligned} P(y|x) &= \frac{P(x|y) P(y)}{P(x)} \\ &= \frac{P(x|y) P(y)}{\int P(x, y) dy} \\ &= \frac{P(x|y) P(y)}{\int P(x|y) P(y) dy} \end{aligned}$$

**Expectation** Expectation tells us the expected or average value of some function  $f(x)$ , taking into account the distribution of  $x$ .

$$\begin{aligned} \mathbf{E}[f(x)] &= \sum_x f(x) P(x) \\ \mathbf{E}[f(x)] &= \int f(x) P(x) dx \end{aligned}$$

Definition in two dimensions:  $\mathbf{E}[f(x, y)] = \iint f(x, y) P(x, y) dx dy$

Function $f(\bullet)$	Expectation
$x$	mean, $\mu_x$
$(x - \mu_x)^2$	variance
$(x - \mu_x)^3$	skew
$(x - \mu_x)^4$	kurtosis
$(x - \mu_x)(x - \mu_y)$	covariance of $x$ and $y$

Besides, Expectation has the following four rules



1. Expected value of a constant is the constant  $\mathbf{E}[\kappa] = \kappa$ .
2. Expected value of constant times function is constant times expected value of function  $\mathbf{E}[kf(x)] = k\mathbf{E}[f(x)]$ .
3. Expectation of sum of functions is sum of expectation of functions  $\mathbf{E}[f(x) + g(y)] = \mathbf{E}[f(x)] + \mathbf{E}[g(y)]$ .
4. Expectation of product of functions in variables  $x$  and  $y$  is product of expectations of functions if  $x$  and  $y$  are independent  $\mathbf{E}[f(x)g(y)] = \mathbf{E}[f(x)]\mathbf{E}[g(y)]$ ,  $x \perp y$ .

### 2.1.2 Common Probability Distributions

**Bernoulli** Bernoulli distribution describes situation where only two possible outcomes  $y = 0/y = 1$  or failure/success.

1.  $P(x) = \mathbf{Bern}_x[\lambda] = \lambda^x(1 - \lambda)^{1-x}$
2. univariate, discrete, binary
3.  $x \in \{0, 1\}; \lambda \in [0, 1]$
4.  $\mathbf{E}[x] = \lambda$ ,  $\mathbf{Var}[x] = \lambda(1 - \lambda)$

**Beta** Beta distribution is the conjugate distribution to Bernoulli distribution.

1.  $P(\lambda) = \mathbf{Beta}_\lambda[\alpha, \beta] = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \lambda^{\alpha-1}(1 - \lambda)^{\beta-1}$
2. univariate, continuous, unbounded
3.  $\lambda \in \mathbb{R}; \alpha \in \mathbb{R}_+, \beta \in \mathbb{R}_+$
4.  $\mathbf{E}[\lambda] = \frac{\alpha}{\alpha+\beta}$ ,  $\mathbf{Var}[\lambda] = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$

**Categorical** Categorical distribution describes situation with  $K$  possible outcomes.

1.  $P(x) = \mathbf{Cat}_x[\boldsymbol{\lambda}]$ ,  $P(x = k) = \lambda_k$ ,  $P(\mathbf{x} = \mathbf{e}_k) = \prod_{j=1}^K \lambda_j^{x_j} = \lambda_k$
2. univariate, discrete, multi-valued
3.  $x \in \{1, 2, \dots, K\}; \lambda_k \in [0, 1]$  where  $\sum_k \lambda_k = 1$
4.  $\mathbf{E}[x_i] = \lambda_i$ ,  $\mathbf{Var}[x_i] = \lambda_i(1 - \lambda_i)$ ,  $\mathbf{Cov}[x_i, x_j] = -\lambda_i\lambda_j$  ( $i \neq j$ )

**Dirichlet** Dirichlet distribution is the conjugate distribution to categorical distribution.

1.  $P(\boldsymbol{\lambda}) = \mathbf{Dir}_\lambda[\boldsymbol{\alpha}] = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \lambda_k^{\alpha_k-1}$
2. multivariate, continuous, bounded, sums to one
3.  $\boldsymbol{\lambda} = [\lambda_1, \lambda_2, \dots, \lambda_K]^\top$ ,  $\lambda_k \in [0, 1]$ ,  $\sum_{k=1}^K \lambda_k = 1$ ;  $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_K]$ ,  $\alpha_k \in \mathbb{R}_+$

$$4. \mathbf{E}[\lambda_i] = \frac{\alpha_i}{\sum_k \alpha_k}, \mathbf{Var}[\lambda_i] = \frac{\alpha_i(\sum_k \alpha_k - \alpha_i)}{(\sum_k \alpha_k)^2(\sum_k \alpha_k + 1)}, \mathbf{Cov}[\lambda_i, \lambda_j] = \frac{-\alpha_i \alpha_j}{(\sum_k \alpha_k)^2(\sum_k \alpha_k + 1)} \quad (i \neq j)$$

**Univariate Normal** Univariate normal distribution describes single continuous variable.

1.  $P(x) = \mathbf{Norm}_x[\mu, \sigma^2] = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$
2. univariate, continuous, unbounded
3.  $x \in \mathbb{R}; \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+$
4.  $\mathbf{E}[x] = \mu, \mathbf{Var}[x] = \sigma^2$

**Normal Inverse Gamma** Normal inverse gamma distribution is a conjugate distribution to univariate normal distribution.

1.  $P(\mu, \sigma^2) = \mathbf{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta] = \frac{\sqrt{\gamma}}{\sqrt{2\pi\sigma^2}} \frac{\beta^\alpha}{\Gamma(\alpha)} \left(\frac{1}{\sigma^2}\right)^{\alpha+1} \exp\left(-\frac{2\beta + \gamma(\delta - \mu)^2}{2\sigma^2}\right)$
2. bivariate, continuous,  $\mu$  unbounded,  $\sigma^2$  bounded below
3.  $\mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+; \alpha \in \mathbb{R}_+, \beta \in \mathbb{R}_+, \gamma \in \mathbb{R}_+, \delta \in \mathbb{R}$
4.  $\mathbf{E}[\mu] = \delta, \mathbf{E}[\sigma^2] = \frac{\beta}{\alpha-1} \quad (\alpha > 1), \mathbf{Var}[\mu] = \frac{\beta}{(\alpha-1)\gamma} \quad (\alpha > 1), \mathbf{Var}[\sigma^2] = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)} \quad (\alpha > 2), \mathbf{Cov}[\mu, \sigma^2] = 0 \quad (\alpha > 1)$

**Multivariate Normal** Multivariate normal distribution describes multiple continuous variables. It takes two parameters: a vector containing mean position  $\boldsymbol{\mu}$ , and a symmetric positive definite covariance matrix  $\boldsymbol{\Sigma}$ .

1.  $P(\mathbf{x}) = \mathbf{Norm}_x[\boldsymbol{\mu}, \boldsymbol{\Sigma}] = \frac{1}{\sqrt{\det(2\pi\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$
2. multivariate, continuous, unbounded
3.  $\mathbf{x} \in \mathbb{R}^K; \boldsymbol{\mu} \in \mathbb{R}^K, \boldsymbol{\Sigma} \in \mathbb{R}^{K \times K}$  (positive semi-definite matrix)
4.  $\mathbf{E}[\mathbf{x}] = \boldsymbol{\mu}, \mathbf{Var}[\mathbf{x}] = \boldsymbol{\Sigma}$

**Normal Inverse Wishart** Normal inverse wishart distribution is a conjugate distribution to multivariate normal distribution.

1.  $P(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathbf{NormInvWis}_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}[\alpha, \boldsymbol{\Psi}, \gamma, \boldsymbol{\delta}]$   
 $= \frac{\gamma^{D/2} |\boldsymbol{\Psi}|^{\alpha/2} |\boldsymbol{\Sigma}|^{-\frac{\alpha+D+2}{2}}}{(2\pi)^{D/2} 2^{(\alpha\boldsymbol{\Sigma})/2} \Gamma_D(\alpha/2)} \exp\left(-\frac{1}{2}(\text{Tr}(\boldsymbol{\Psi}\boldsymbol{\Sigma}^{-1}) + \gamma(\boldsymbol{\mu} - \boldsymbol{\delta})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \boldsymbol{\delta}))\right)$
2. multivariate,  $\boldsymbol{\mu}$  unbounded,  $\boldsymbol{\Sigma}$  square, positive definite
3.  $\boldsymbol{\mu} \in \mathbb{R}^K, \boldsymbol{\Sigma} \in \mathbb{R}^{K \times K}; \alpha \in \mathbb{R}_{>D-1}, \boldsymbol{\Psi} \in \mathbb{R}^{K \times K}, \gamma \in \mathbb{R}_+, \boldsymbol{\delta} \in \mathbb{R}^K$

## 2.2 Linear Algebra

### 2.2.1 Vectors

#### Vectors Addition

$$\mathbf{v} + \mathbf{w} = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} + \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} v_1 + w_1 \\ v_2 + w_2 \\ \vdots \\ v_n + w_n \end{pmatrix}$$

#### Vectors Scaling

$$a\mathbf{v} = a \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} = \begin{pmatrix} av_1 \\ av_2 \\ \vdots \\ av_n \end{pmatrix}$$

#### Rules for Vectors Addition and Scaling

1.  $\mathbf{u} + (\mathbf{v} + \mathbf{w}) = (\mathbf{u} + \mathbf{v}) + \mathbf{w}$
2.  $\mathbf{v} + \mathbf{w} = \mathbf{w} + \mathbf{v}$
3. There is a vector  $\mathbf{0}$  such that  $\mathbf{0} + \mathbf{v} = \mathbf{v}$  for all  $\mathbf{v}$
4. For every vector  $\mathbf{v}$  there is a vector  $-\mathbf{v}$  so that  $\mathbf{v} + (-\mathbf{v}) = \mathbf{0}$
5.  $a(b\mathbf{v}) = (ab)\mathbf{v}$
6.  $1\mathbf{v} = \mathbf{v}$
7.  $a(\mathbf{v} + \mathbf{w}) = a\mathbf{v} + a\mathbf{w}$
8.  $(a + b)\mathbf{v} = a\mathbf{v} + b\mathbf{v}$

#### Linear Combination & Span

Linear combination (e.g. in 2D space):

$$a\mathbf{v} + b\mathbf{w}$$

The span of  $\mathbf{v}$  and  $\mathbf{w}$  is the set of all their linear combinations.

## Representation of Basis

In Euclidean:

$$\begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} = v_1 \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} + v_2 \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix} + \cdots + v_n \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$$

We can write this as

$$\mathbf{v} = v_1 \mathbf{e}^1 + v_2 \mathbf{e}^2 + \cdots + v_n \mathbf{e}^n$$

In different basis, we choose other basis vector and then write the same vector

$$\mathbf{v} = w_1 \mathbf{b}^1 + w_2 \mathbf{b}^2 + \cdots + w_n \mathbf{b}^n$$

If these basis vectors are orthonormal,  $w_i = \mathbf{v}^\top \mathbf{b}^i$

## Linear Dependence

1. Linearly dependent: A set of vectors  $\mathbf{v}^1, \dots, \mathbf{v}^n$  is linearly dependent if there exists a vector  $\mathbf{v}^j$  that can be expressed as a linear combination of the other vectors. (The vector  $\mathbf{v}^j$  is already located in the span of other vectors)
2. Linearly Independent: Each vector really does add another dimension to the span. And the only solution to  $\sum_{i=1}^n a_i \mathbf{v}^i = \mathbf{0}$  is for all  $a_i = 0, i = 1, \dots, n$ .

## Dot Products

$$\mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^n v_i w_i = \mathbf{v}^\top \mathbf{w}$$

The length of a vector is denoted as  $\|\mathbf{v}\|$ , the squared length is given by

$$\|\mathbf{v}\|^2 = \mathbf{v}^\top \mathbf{v} = \mathbf{v}^2 = v_1^2 + v_2^2 + \cdots + v_n^2$$

A natural geometric interpretation of dot products is

$$\mathbf{v} \cdot \mathbf{w} = \|\mathbf{v}\| \|\mathbf{w}\| \cos \theta$$

where  $\theta$  is the angle between two vectors.

### **2.2.2 Matrices**

## **2.3 Calculus**

### **2.3.1 Differentiation & Integration**

### **2.3.2 Multivariate Calculus**

## **2.4 Informatics**

### **2.4.1 Entropy**

## **2.5 Optimization**

### **2.5.1 One-dimensional Minimization**

### **2.5.2 Gradient Descent**

### **2.5.3 Quadratic Functions**

### **2.5.4 General Functions**

### **2.5.5 Optimization with Constraints**

# Chapter 3

## Machine Learning Basics

### 3.1 Regularization

#### 3.1.1 Under-fitting & Over-fitting

**Under-fitting** If  $N > D$  (e.g. 30 data points, 2 dimensions) we have more equations than unknowns: over-determined system. Input-output relations can only hold approximately.

**Over-fitting** If  $N < D$  (e.g. 30 points, 15265 dimensions) we have more unknowns than equations: under-determined system. Input-output equations hold exactly, but we are simply memorizing data.

#### 3.1.2 Bias & Variance

**High Bias & Low Variance** A rigid model's (low complexity) performance is more predictable in the test set but the model may not be good even on the training set.

**Low Bias & High Variance** A flexible model (high complexity) approximates the target function well in the training set but can "overtrain" and have poor performance on the test set.

#### 3.1.3 Vector Norm

**L1, ("Manhattan") norm**  $\|\mathbf{w}\|_1 = \sum_{d=1}^D |w_d|$

**L2, ("Euclidean") norm**  $\|\mathbf{w}\|_2 = \sqrt{\sum_{d=1}^D w_d^2} = \sqrt{\langle \mathbf{w}, \mathbf{w} \rangle} = \sqrt{\mathbf{w}^\top \mathbf{w}}$

**Lp norm,  $p > 1$**   $\|\mathbf{w}\|_p = \left( \sum_{d=1}^D w_d^p \right)^{1/p}$

### 3.1.4 Penalize Complexity

In linear regression, the residual vector is  $\epsilon = \mathbf{y} - \Psi\mathbf{w}$ . The loss function is  $L(\mathbf{w}) = \epsilon^\top \epsilon$ . We add a complexity term  $R(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \mathbf{w}^\top \mathbf{w}$  to the loss function. Hence, the original loss function becomes  $L(\mathbf{w}) = \epsilon^\top \epsilon + \lambda \mathbf{w}^\top \mathbf{w}$ .

Without regularization, the loss function is  $L(\mathbf{w}) = \epsilon^\top \epsilon$ . Let  $\nabla L(\mathbf{w}^*) = 0$ , we have  $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ .

With L2-regularization, the loss function is  $L(\mathbf{w}) = \epsilon^\top \epsilon + \lambda \mathbf{w}^\top \mathbf{w}$ . Let  $\nabla L(\mathbf{w}^*) = 0$ , we have  $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$ . The additional  $\lambda \mathbf{I}$  makes the data matrix more robust to calculate inversion.

## 3.2 Cross-Validation

We can select hyperparameters with (cross-)validation. Cross-validation excludes part of the training data from parameter estimation, and use them only to predict the test error.

K-fold cross validation: split data set into K folds and each time train on (K-1) folds and valid on the remaining fold until all folds have been used as validation fold. The cross-validation error is the average of K validation errors. We pick hyperparameters that minimize cross-validation error.

## 3.3 Bayesian Learning

### 3.3.1 Bayes' Rule Terminology

Bayes' Rule:

$$P(y|x) = \frac{P(x|y) P(y)}{\int P(x|y) P(y) dy}$$

**Prior**  $P(y)$  what we know about  $y$  before seeing  $x$ . In parameters learning we choose prior that is conjugate to likelihood.

**Likelihood**  $P(x|y)$  propensity for observing a certain value of  $x$  given a certain value of  $y$ .

**Posterior**  $P(y|x)$  what we know about  $y$  after seeing  $x$ . Posterior must have same form as conjugate prior distribution.

**Evidence**  $\int P(x|y) P(y) dy$  a constant to ensure that the LHS is a valid distribution. Posterior must be a distribution which implies that evidence equals to a constant  $\kappa$  from conjugate relation.

### 3.3.2 Maximum Likelihood

**Fitting** As the name suggests we find the parameters under which the data  $\mathbf{x}_{1...I}$  are most likely. Here, we have assumed that data was independent.

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} P(\mathbf{x}_{1...I}|\theta) \\ &= \operatorname{argmax}_{\theta} \prod_{i=1}^I P(\mathbf{x}_i|\theta)\end{aligned}$$

**Predictive Density** Evaluate new data point  $\mathbf{x}^*$  under probability distribution  $P(\mathbf{x}^*|\hat{\theta})$  with best parameters.

### 3.3.3 Maximum a Posterior (MAP)

**Fitting** As the name suggests we find the parameters which maximize the posterior probability  $P(\theta|\mathbf{x}_{1...I})$ . Again we have assumed that data was independent.

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} P(\theta|\mathbf{x}_{1...I}) \\ &= \operatorname{argmax}_{\theta} \frac{P(\mathbf{x}_{1...I}|\theta) P(\theta)}{P(\mathbf{x}_{1...I})} \\ &= \operatorname{argmax}_{\theta} \frac{\prod_{i=1}^I P(\mathbf{x}_i|\theta) P(\theta)}{P(\mathbf{x}_{1...I})}\end{aligned}$$

Since the denominator does not depend on the parameters we can instead maximize

$$\hat{\theta} = \operatorname{argmax}_{\theta} \prod_{i=1}^I P(\mathbf{x}_i|\theta) P(\theta)$$

**Predictive Density** Evaluate new data point  $\mathbf{x}^*$  under probability distribution with MAP parameters  $P(\mathbf{x}^*|\hat{\theta})$

### 3.3.4 Bayesian Approach

**Fitting** Compute the posterior distribution over possible parameter values using Bayes' rule. Principle: There are many values that could have explained the data. Instead of picking one set of parameters, try to capture all of the possibilities.

$$P(\theta|\mathbf{x}_{1...I}) = \frac{\prod_{i=1}^I P(\mathbf{x}_i|\theta) P(\theta)}{P(\mathbf{x}_{1...I})}$$



**Predictive Density** (a) Each possible parameter value makes a prediction. (b) Some parameters more probable than others.

$$P(\mathbf{x}^*|\mathbf{x}_{1...I}) = \int P(\mathbf{x}^*|\boldsymbol{\theta}) P(\boldsymbol{\theta}|\mathbf{x}_{1...I}) d\boldsymbol{\theta}$$

Make a prediction that is an infinite weighted sum (integral) of the predictions for each parameter value ( $P(\mathbf{x}^*|\boldsymbol{\theta})$ ), where weights are the probabilities ( $P(\boldsymbol{\theta}|\mathbf{x}_{1...I})$ ).

### 3.3.5 Example: Univariate Normal Distribution

#### Maximum Likelihood

Likelihood given by normal distribution pdf:

$$P(x|\mu, \sigma^2) = \mathbf{Norm}_x[\mu, \sigma^2] = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Apply maximum likelihood:

$$\begin{aligned} \hat{\mu}, \hat{\sigma}^2 &= \operatorname{argmax}_{\mu, \sigma^2} P(\mathbf{x}_{1...I}|\mu, \sigma^2) \\ &= \operatorname{argmax}_{\mu, \sigma^2} \prod_{i=1}^I P(x_i|\mu, \sigma^2) \\ &= \operatorname{argmax}_{\mu, \sigma^2} \prod_{i=1}^I \mathbf{Norm}_{x_i}[\mu, \sigma^2] \\ &= \operatorname{argmax}_{\mu, \sigma^2} \sum_{i=1}^I \log \mathbf{Norm}_{x_i}[\mu, \sigma^2] \\ &= \operatorname{argmax}_{\mu, \sigma^2} \left( -\frac{I}{2} \log 2\pi - \frac{I}{2} \log \sigma^2 - \frac{1}{2} \sum_{i=1}^I \frac{(x_i - \mu)^2}{\sigma^2} \right) \end{aligned}$$

Let  $\nabla L(\hat{\mu}, \hat{\sigma}^2) = 0$ , we have the solution:

$$\begin{aligned} \hat{\mu} &= \frac{\sum_{i=1}^I x_i}{I} \\ \hat{\sigma}^2 &= \sum_{i=1}^I \frac{(x_i - \hat{\mu})^2}{I} \end{aligned}$$

#### Maximum a Posterior

Likelihood given by normal distribution pdf:

$$P(x|\mu, \sigma^2) = \mathbf{Norm}_x[\mu, \sigma^2] = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Prior given by normal inverse gamma distribution pdf:

$$P(\mu, \sigma^2) = \mathbf{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta] = \frac{\sqrt{\gamma}}{\sqrt{2\pi\sigma^2}} \frac{\beta^\alpha}{\Gamma(\alpha)} \left(\frac{1}{\sigma^2}\right)^{\alpha+1} \exp\left(-\frac{2\beta + \gamma(\delta - \mu)^2}{2\sigma^2}\right)$$

Apply maximum a posterior:

$$\begin{aligned} \hat{\mu}, \hat{\sigma}^2 &= \underset{\mu, \sigma^2}{\operatorname{argmax}} \prod_{i=1}^I P(x_i | \mu, \sigma^2) P(\mu, \sigma^2) \\ &= \underset{\mu, \sigma^2}{\operatorname{argmax}} \prod_{i=1}^I \mathbf{Norm}_{x_i}[\mu, \sigma^2] \mathbf{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta] \\ &= \underset{\mu, \sigma^2}{\operatorname{argmax}} \left( \sum_{i=1}^I \log \mathbf{Norm}_{x_i}[\mu, \sigma^2] + \log \mathbf{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta] \right) \end{aligned}$$

Let  $\nabla L(\hat{\mu}, \hat{\sigma}^2) = 0$ , we have the solution:

$$\begin{aligned} \hat{\mu} &= \frac{\sum_{i=1}^I x_i + \gamma\delta}{I + \gamma} \\ \hat{\sigma}^2 &= \frac{\sum_{i=1}^I (x_i - \mu)^2 + 2\beta + \gamma(\delta - \mu)^2}{I + 3 + 2\alpha} \end{aligned}$$

## Bayesian Approach

Compute the posterior distribution using Bayes' rule:

$$\begin{aligned} P(\mu, \sigma^2 | x_{1...I}) &= \frac{\prod_{i=1}^I P(x_i | \mu, \sigma^2) P(\mu, \sigma^2)}{P(x_{1...I})} \\ &= \frac{\prod_{i=1}^I \mathbf{Norm}_{x_i}[\mu, \sigma^2] \mathbf{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta]}{P(x_{1...I})} \\ &= \frac{\kappa(\alpha, \beta, \gamma, \delta, x_{1...I}) \mathbf{NormInvGam}_{\mu, \sigma^2}[\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma}, \tilde{\delta}]}{P(x_{1...I})} \\ &= \mathbf{NormInvGam}_{\mu, \sigma^2}[\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma}, \tilde{\delta}] \end{aligned}$$

where

$$\begin{aligned} \tilde{\alpha} &= \alpha + \frac{I}{2} \\ \tilde{\beta} &= \frac{\sum_i x_i^2}{2} + \beta + \frac{\gamma\delta^2}{2} - \frac{(\gamma\delta + \sum_i x_i)^2}{2(\gamma + I)} \\ \tilde{\gamma} &= \gamma + I \\ \tilde{\delta} &= \frac{\gamma\delta + \sum_i x_i}{\gamma + I} \end{aligned}$$

Take weighted sum of predictions from different parameter values:

$$\begin{aligned}
P(x^*|x_{1...I}) &= \iint P(x^*|\mu, \sigma^2) P(\mu, \sigma^2|x_{1...I}) d\mu d\sigma \\
&= \iint \mathbf{Norm}_{x^*}[\mu, \sigma^2] \mathbf{NormInvGam}_{\mu, \sigma^2}[\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma}, \tilde{\delta}] d\mu d\sigma \\
&= \iint \kappa(\alpha, \beta, \gamma, \delta, x_{1...I}) \mathbf{NormInvGam}_{\mu, \sigma^2}[\check{\alpha}, \check{\beta}, \check{\gamma}, \check{\delta}] d\mu d\sigma \\
&= \kappa(\alpha, \beta, \gamma, \delta, x_{1...I}) \iint \mathbf{NormInvGam}_{\mu, \sigma^2}[\check{\alpha}, \check{\beta}, \check{\gamma}, \check{\delta}] d\mu d\sigma \\
&= \kappa(\alpha, \beta, \gamma, \delta, x_{1...I}) \\
&= \frac{1}{\sqrt{2\pi}} \frac{\sqrt{\tilde{\gamma}} \tilde{\beta}^{\tilde{\alpha}} \Gamma(\check{\alpha})}{\sqrt{\check{\gamma}} \check{\beta}^{\check{\alpha}} \Gamma(\check{\alpha})}
\end{aligned}$$

where

$$\begin{aligned}
\check{\alpha} &= \tilde{\alpha} + \frac{1}{2} \\
\check{\beta} &= \frac{x^{*2}}{2} + \tilde{\beta} + \frac{\tilde{\gamma}^{\tilde{\delta}^2}}{2} - \frac{(\tilde{\gamma}^{\tilde{\delta}} + x^*)^2}{2(\tilde{\gamma} + 1)} \\
\check{\gamma} &= \tilde{\gamma} + 1
\end{aligned}$$

### 3.3.6 Example: Categorical Distribution

#### Maximum Likelihood

Likelihood given by categorical distribution pdf:

$$P(x|\boldsymbol{\lambda}) = \mathbf{Cat}_x[\boldsymbol{\lambda}] = \prod_{j=1}^K \lambda_j^{x_j} = \lambda_k$$

Apply maximum likelihood:

$$\begin{aligned}
\hat{\boldsymbol{\lambda}} &= \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \prod_{i=1}^I P(x_i|\boldsymbol{\lambda}) & s.t. \sum_k \lambda_k &= 1 \\
&= \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \prod_{i=1}^I \mathbf{Cat}_{x_i}[\boldsymbol{\lambda}] & s.t. \sum_k \lambda_k &= 1 \\
&= \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \prod_{k=1}^K \lambda_k^{N_k} & s.t. \sum_k \lambda_k &= 1 \\
&= \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \sum_{k=1}^K N_k \log \lambda_k & s.t. \sum_k \lambda_k &= 1
\end{aligned}$$

Here,  $N_k$  represents the number of times the data is classified in class  $k$ . As before, we will instead optimize log probability. Since there is a constraint *s.t.*  $\sum_k \lambda_k = 1$ , we use Lagrange multiplier to reconstruct the loss function.

$$L(\boldsymbol{\lambda}) = \sum_{k=1}^K N_k \log \lambda_k + v \left( \sum_{k=1}^K \lambda_k - 1 \right)$$

Let  $\nabla L(\boldsymbol{\lambda}, v) = 0$ , we have the solution:

$$\hat{\lambda}_k = \frac{N_k}{\sum_{m=1}^K N_m}$$

### Maximum a Posterior

Likelihood given by categorical distribution pdf:

$$P(x|\boldsymbol{\lambda}) = \mathbf{Cat}_x[\boldsymbol{\lambda}] = \prod_{j=1}^K \lambda_j^{x_j} = \lambda_k$$

Prior given by Dirichlet distribution pdf:

$$P(\boldsymbol{\lambda}) = \mathbf{Dir}_{\boldsymbol{\lambda}}[\boldsymbol{\alpha}] = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \lambda_k^{\alpha_k - 1}$$

Apply maximum a posterior:

$$\begin{aligned} \hat{\boldsymbol{\lambda}} &= \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \prod_{i=1}^I P(x_i|\boldsymbol{\lambda}) P(\boldsymbol{\lambda}) & \text{s.t. } \sum_k \lambda_k &= 1 \\ &= \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \mathbf{Cat}_{x_i}[\boldsymbol{\lambda}] \mathbf{Dir}_{\boldsymbol{\lambda}}[\boldsymbol{\alpha}] & \text{s.t. } \sum_k \lambda_k &= 1 \\ &= \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \prod_{k=1}^K \lambda_k^{N_k} \prod_{k=1}^K \lambda_k^{\alpha_k - 1} & \text{s.t. } \sum_k \lambda_k &= 1 \\ &= \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \prod_{k=1}^K \lambda_k^{N_k + \alpha_k - 1} & \text{s.t. } \sum_k \lambda_k &= 1 \\ &= \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \sum_{k=1}^K (N_k + \alpha_k - 1) \log \lambda_k & \text{s.t. } \sum_k \lambda_k &= 1 \end{aligned}$$

The loss function is very similar to maximum likelihood (same when the prior is uniform, i.e.  $\alpha_{1\dots k} = 1$ ). Take derivative with Lagrange multiplier, we have the solution:

$$\hat{\lambda}_k = \frac{N_k + \alpha_k - 1}{\sum_{m=1}^K (N_m + \alpha_m - 1)}$$

## Bayesian Approach

Compute the posterior distribution using Bayes' rule:

$$\begin{aligned}
 P(\boldsymbol{\lambda}|x_{1...I}) &= \frac{\prod_{i=1}^I P(x_i|\boldsymbol{\lambda}) P(\boldsymbol{\lambda})}{P(x_{1...I})} \\
 &= \frac{\prod_{i=1}^I \text{Cat}_{x_i}[\boldsymbol{\lambda}] \text{Dir}_{\boldsymbol{\lambda}}[\boldsymbol{\alpha}]}{P(x_{1...I})} \\
 &= \frac{\kappa(\boldsymbol{\alpha}, x_{1...I}) \text{Dir}_{\boldsymbol{\lambda}}[\tilde{\boldsymbol{\alpha}}]}{P(x_{1...I})} \\
 &= \text{Dir}_{\boldsymbol{\lambda}}[\tilde{\boldsymbol{\alpha}}]
 \end{aligned}$$

Compute predictive distribution:

$$\begin{aligned}
 P(x^*|x_{1...I}) &= \int P(x^*|\boldsymbol{\lambda}) P(\boldsymbol{\lambda}|x_{1...I}) d\boldsymbol{\lambda} \\
 &= \int \text{Cat}_{x^*}[\boldsymbol{\lambda}] \text{Dir}_{\boldsymbol{\lambda}}[\tilde{\boldsymbol{\alpha}}] d\boldsymbol{\lambda} \\
 &= \int \kappa(x^*, \tilde{\boldsymbol{\alpha}}) \text{Dir}_{\boldsymbol{\lambda}}[\tilde{\boldsymbol{\alpha}}] d\boldsymbol{\lambda} \\
 &= \kappa(x^*, \boldsymbol{\alpha})
 \end{aligned}$$

## 3.4 Machine Learning Models

### 3.4.1 Learning and Inference

In real world problems, we usually have two tasks:

1. Observe measured data,  $\mathbf{x}$
2. Draw inferences from it about world,  $\mathbf{w}$

and

1. When the world state  $\mathbf{w}$  is *continuous*, we'll call this *regression*.
2. When the world state  $\mathbf{w}$  is *discrete*, we'll call this *classification*.

We want take observations  $\mathbf{x}$ , and return probability distribution  $P(\mathbf{w}|\mathbf{x})$  over possible worlds compatible with data. To solve this, we need

1. A *model* that mathematically relates the visual data  $\mathbf{x}$  to the world state  $\mathbf{w}$ . Model specifies family of relationships, particular relationship depends on parameter  $\theta$ .
2. A *learning algorithm* fits parameters  $\theta$  from paired training examples  $\mathbf{x}_i, \mathbf{w}_i$ .
3. An *inference algorithm* uses model to return  $P(\mathbf{w}|\mathbf{x})$  given new observation data  $\mathbf{x}$ .

### 3.4.2 Three Types of Model

We have three types of model:

1. Model contingency of the world on the data  $P(w|x)$ . (Discriminative Model)
2. Model joint occurrence of world and data  $P(x, w)$ . (Generative Model)
3. Model contingency of data on world  $P(x|w)$ . (Generative Model)

Within the three models, type 1 is called *Discriminative Model*. Type 2 and 3 are called *Generative Model*.

#### Model $P(w|x)$ - Discriminative

1.  $P(w|x, \theta) = \mathbf{Distrib}_w[f(x, \theta)]$
2. How to model: (a) Choose an appropriate form for  $P(w)$ . (b) Make parameters a function of  $x$ . (c) Function takes parameters  $\theta$  that define its shape.
3. Learning algorithm: Learn parameters  $\theta$  from training data  $x, w$ .
4. Inference algorithm: Just evaluate  $P(w|x)$

#### Model $P(x, w)$ - Generative

1.  $P(z|\theta) = \mathbf{Distrib}_z[\theta]$
2. How to model: (a) Concatenate  $x$  and  $w$  to make  $z = [x^\top, w^\top]^\top$ . (b) Model the pdf of  $z$ . (c) pdf takes parameters  $\theta$  that define its shape.
3. Learning algorithm: Learn parameters  $\theta$  from training data  $x, w$ .
4. Inference algorithm: Compute  $P(w|x)$  using Bayes' rule  $P(w|x) = \frac{P(x, w)}{P(x)} = \frac{P(x, w)}{\int P(x, w) dw}$ .

#### Model $P(x|w)$ - Generative

1.  $P(x|w, \theta) = \mathbf{Distrib}_x[f(w, \theta)]$
2. How to model: (a) Choose an appropriate form for  $P(x)$ . (b) Make parameters a function of  $w$ . (c) Function takes parameters  $\theta$  that define its shape.
3. Learning algorithm: Learn parameters  $\theta$  from training data  $x, w$ .
4. Define prior  $P(w)$  and then compute  $P(w|x)$  using Bayes' rule  $P(w|x) = \frac{P(x|w)P(w)}{\int P(x|w)P(w)dw}$ .

### 3.4.3 Example: Regression

Consider a simple case:

1. We make a univariate continuous measurement  $x$ .
2. Use this to predict a univariate continuous state  $w$ .

#### Model $P(w|x)$ - Discriminative

1.  $P(w|x, \theta) = \text{Norm}_w[\phi_0 + \phi_1 x, \sigma^2], \theta = \{\phi_0, \phi_1, \sigma^2\}$
2. How to model: (a) Choose normal distribution over  $w$ . (b) Make mean  $\mu$  linear function of  $x$  (variance constant). (c) Parameters are  $\phi_0$  (y-offset),  $\phi_1$  (slope),  $\sigma^2$  (variance). This model is called *linear regression*.
3. Learning algorithm: Learn  $\theta$  from training data  $x, w$ . e.g. MAP:

$$\begin{aligned}\hat{\theta} &= \underset{\theta}{\operatorname{argmax}} P(\theta | w_{1...I}, x_{1...I}) \\ &= \underset{\theta}{\operatorname{argmax}} P(w_{1...I} | x_{1...I}, \theta) P(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^I P(w_i | x_i, \theta) P(\theta)\end{aligned}$$

4. Inference algorithm: Just evaluate  $P(w|x)$  for new data  $x$ .

#### Model $P(x, w)$ - Generative

1.  $P(x, w | \theta) = \text{Norm}_{x,w}[\mu, \Sigma], \theta = \{\mu, \Sigma\}$
2. How to model: (a) Concatenate  $x$  and  $w$  to make  $z = [x^\top, w^\top]^\top$ . (b) Model the pdf of  $z$  as normal distribution. (c) pdf makes parameters  $\mu$  and  $\Sigma$  that define its shape.
3. Learning algorithm: Learn parameters  $\theta$  from training data  $x, w$ .
4. Inference algorithm: Compute  $P(w|x)$  using Bayes' rule  $P(w|x) = \frac{P(x,w)}{P(x)} = \frac{P(x,w)}{\int P(x,w) dw}$ .

#### Model $P(x|w)$ - Generative

1.  $P(x|w, \theta) = \text{Norm}_x[\phi_0 + \phi_1 w, \sigma^2], \theta = \{\phi_0, \phi_1, \sigma^2\}$
2. How to model: (a) Choose normal distribution over  $x$ . (b) Make mean  $\mu$  linear function of  $w$  (variance constant). (c) Parameters are  $\phi_0, \phi_1, \sigma^2$ .

3. Learning algorithm: Learn  $\theta$  from training data  $x, w$ . e.g. MAP
4. Inference algorithm: Compute  $P(w|x)$  using Bayes' rule  $P(w|x) = \frac{P(x,w)}{P(x)} = \frac{P(x,w)}{\int P(x,w)dw}$ .

### 3.4.4 Example: Classification

Consider a simple case:

1. We make a univariate continuous measurement  $x$ .
2. Use this to predict a discrete binary world  $w \in \{0, 1\}$ .

#### Model $P(w|x)$ - Discriminative

1.  $P(w|x, \theta) = \text{Bern}_w[\sigma(\phi_0 + \phi_1 x)], \theta = \phi_0, \phi_1$
2. How to model: (a) Choose Bernoulli distribution for  $P(w)$ . (b) Make parameters a sigmoid-activated function of  $x$ . (c) Function takes parameters  $\phi_0$  and  $\phi_1$ . This model is called *logistic regression*.
3. Learning algorithm: Learning by standard methods, e.g. ML, MAP, Bayesian Approach.
4. Inference algorithm: Just evaluate  $P(w|x)$ .

#### Model $P(x, w)$ - Generative

Can't build this mode very easily:

1. Concatenate continuous vector  $x$  and discrete  $w$  to make  $z$ .
2. No obvious probability distribution to model joint probability of discrete and continuous.

#### Model $P(x|w)$ - Generative

1.  $P(x|w, \theta) = \text{Norm}_x[\mu_w, \sigma_w^2], \theta = \{\mu_0, \mu_1, \sigma_0^2, \sigma_1^2\}$
2. How to model: (a) Choose a Normal distribution for  $P(x)$ . (b) Make parameters a function of discrete binary  $w$ . (c) Function takes parameters  $\mu_0, \mu_1, \sigma_0^2, \sigma_1^2$  that define its shape.
3. Learning algorithm: Learning by standard methods, e.g. ML, MAP, Bayesian Approach.
4. Define prior  $P(w)$  and then compute  $P(w|x)$  using Bayes' rule  $P(w|x) = \frac{P(x|w)P(w)}{\int P(x|w)P(w)dw}$ .



## 3.5 Overview of Common Algorithms

Properties of common machine learning methods:

Method <sup>1</sup>	Problem <sup>2</sup>	Model <sup>3</sup>	Learning <sup>4</sup>	Loss Function	Algorithm <sup>5</sup>
Perceptron	BC	D			SGD
K-NN	MC, R	D			
Naive Bayes	MC	G	ML, MAP	$-\log P(w x)$	Bayes, EM
Decision Tree	MC, R	D	NML	$-\log P(w x)$	
LR	MC	D	ML, NML	$\log(1 + \exp(-wf(x)))$	SGD, QN
SVM	BC	D		$[1 - wf(x)]_+$	SMO
Boosting	BC	D		$\exp(-wf(x))$	
EM			ML, MAP	$-\log P(w x)$	Iteration
HMM	T	G	ML, MAP	$-\log P(w x)$	Bayes, EM
CRF	T	D	ML, NML	$-\log P(w x)$	SGD, QN

<sup>1</sup> K-NN=K-Nearest Neighbors, LR=Logistic Regression, SVM=Support Vector Machine, HMM=Hidden Markov Model, CRF=Conditional Random Field

<sup>2</sup> BC=Binary Classification, MC=Multi-class Classification, R=Regression, T=Tagging

<sup>3</sup> D=Discriminative Model, G=Generative Model

<sup>4</sup> ML=Maximum Likelihood, NML=Normalized ML MAP=Maximum a Posterior

<sup>5</sup> SGD=Stochastic Gradient Descent, QN=Quasi-Newton

# Chapter 4

## Linear Regression

### 4.1 Basic Model

1. Discriminative, Regression

$$2. P(w_i|\mathbf{X}_i, \boldsymbol{\theta}) = \text{Norm}_{w_i}[\phi_0 + \boldsymbol{\phi}^\top \mathbf{X}_i, \sigma^2]$$

$$3. \text{ (Neater Notation) } \mathbf{X}_i \leftarrow [1 \quad \mathbf{X}_i^\top]^\top, \boldsymbol{\phi} \leftarrow [\phi_0 \quad \boldsymbol{\phi}^\top]^\top$$
$$P(w_i|\mathbf{X}_i, \boldsymbol{\theta}) = \text{Norm}_{w_i}[\boldsymbol{\phi}^\top \mathbf{X}_i, \sigma^2]$$

$$4. \text{ (Combining Equations) } \mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_I]$$
$$P(\mathbf{w}|\mathbf{X}, \boldsymbol{\theta}) = \text{Norm}_{\mathbf{w}}[\mathbf{X}^\top \boldsymbol{\phi}, \sigma^2 \mathbf{I}]$$

5. Learning with Maximum Likelihood:  $\hat{\boldsymbol{\theta}} = \text{argmax}_{\boldsymbol{\theta}} P(\mathbf{w}|\mathbf{X}, \boldsymbol{\theta}) = \text{argmax}_{\boldsymbol{\theta}} \log P(\mathbf{w}|\mathbf{X}, \boldsymbol{\theta})$ ,  
result:

$$\hat{\boldsymbol{\phi}} = (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{w}$$
$$\hat{\sigma}^2 = \frac{(\mathbf{w} - \mathbf{X}^\top \hat{\boldsymbol{\phi}})^\top (\mathbf{w} - \mathbf{X}^\top \hat{\boldsymbol{\phi}})}{\mathbf{I}}$$

### 4.2 Bayesian Regression

**Parameter  $\boldsymbol{\phi}$**

**Likelihood**  $P(\mathbf{w}|\mathbf{X}, \boldsymbol{\theta}) = \text{Norm}_{\mathbf{w}}[\mathbf{X}^\top \boldsymbol{\phi}, \sigma^2 \mathbf{I}]$

**Prior**  $P(\boldsymbol{\phi}) = \text{Norm}_{\boldsymbol{\phi}}[\mathbf{0}, \sigma_p^2 \mathbf{I}]$

**Posterior**

$$\begin{aligned}
P(\phi|\mathbf{X}, \mathbf{w}) &= \frac{P(\mathbf{w}|\mathbf{X}, \phi) P(\phi|\mathbf{X})}{P(\mathbf{w}|\mathbf{X})} \\
&= \text{Norm}_{\phi} \left[ \frac{1}{\sigma^2} \mathbf{A}^{-1} \mathbf{X} \mathbf{w}, \mathbf{A}^{-1} \right]
\end{aligned}$$

where

$$\begin{aligned}
\mathbf{A} &= \frac{1}{\sigma^2} \mathbf{X} \mathbf{X}^{\top} + \frac{1}{\sigma_p^2} \mathbf{I} \\
\mathbf{A}^{-1} &= \sigma_p^2 \mathbf{I}_D - \sigma_p^2 \mathbf{X} \left( \mathbf{X}^{\top} \mathbf{X} + \frac{\sigma^2}{\sigma_p^2} \mathbf{I}_I \right)^{-1} \mathbf{X}^{\top}
\end{aligned}$$

**Inference (Bayesian Approach)**

$$\begin{aligned}
P(w^*|\mathbf{x}^*, \mathbf{X}, \mathbf{w}) &= \int P(w^*|\mathbf{x}^*, \phi) P(\phi|\mathbf{X}, \mathbf{w}) d\phi \\
&= \int \text{Norm}_{w^*}[\phi^{\top} \mathbf{x}^*, \sigma^2] \text{Norm}_{\phi} \left[ \frac{1}{\sigma^2} \mathbf{A}^{-1} \mathbf{X} \mathbf{w}, \mathbf{A}^{-1} \right] d\phi \\
&= \text{Norm}_{w^*} \left[ \frac{1}{\sigma^2} \mathbf{x}^{*\top} \mathbf{A}^{-1} \mathbf{X} \mathbf{w}, \mathbf{x}^{*\top} \mathbf{A}^{-1} \mathbf{x}^* + \sigma^2 \right]
\end{aligned}$$

**Parameter  $\sigma^2$** 

$$\begin{aligned}
P(\mathbf{w}|\mathbf{X}, \sigma^2) &= \int P(\mathbf{w}|\mathbf{X}, \phi, \sigma^2) P(\phi) d\phi \\
&= \int \text{Norm}_{\mathbf{w}}[\mathbf{X}^{\top} \phi, \sigma^2 \mathbf{I}] \text{Norm}_{\phi}[\mathbf{0}, \sigma_p^2 \mathbf{I}] d\phi \\
&= \text{Norm}_{\mathbf{w}}[\mathbf{0}, \sigma_p^2 \mathbf{X}^{\top} \mathbf{X} + \sigma^2 \mathbf{I}]
\end{aligned}$$

**4.3 Non-linear Regression****4.4 Kernel Trick & Gaussian Processes****4.5 Sparse Linear Regression****4.6 Dual Linear Regression****4.7 Relevance Vector Regression**

# Chapter 5

## Logistic Regression

### 5.1 Logistic Regression

### 5.2 Non-linear Logistic Regression

### 5.3 Kernel Trick & Gaussian Process Classification

### 5.4 Multi-class Classification

# Chapter 6

## Support Vector Machines

6.1 Geometric Margins

6.2 Primal & Dual Problems

6.3 Support Vectors

6.4 Slack Variables

6.5 Hinge Loss

6.6 Non-linear SVMs

6.7 Kernel Trick

# Chapter 7

## EM Algorithm

7.1 Expectation Maximization

7.2 Example: Mixture of Gaussians

7.3 Example: t-distributions

7.4 Example: Factor Analysis

# Chapter 8

## Boosting

### 8.1 Ensemble Methods

### 8.2 Bagging

### 8.3 Boosting

### 8.4 Adaboost

# Chapter 9

## Decision Tree & Random Forest

### 9.1 CART

### 9.2 ID3

### 9.3 C4.5

### 9.4 Random Forest



# Chapter 10

## Graphical Models & Markov Network

### 10.1 Graph Definitions

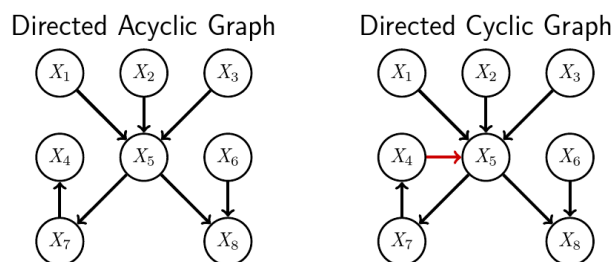
#### 10.1.1 Graph

**Graph** A graph consists of nodes (vertices) and undirected or directed links (edges) between nodes.

**Path** A path from  $X_i$  to  $X_j$  is a sequence of connected nodes starting at  $X_i$  and ending at  $X_j$ .

#### 10.1.2 Directed Graph

**Directed Graphs** Graphs that all the edges are directed.



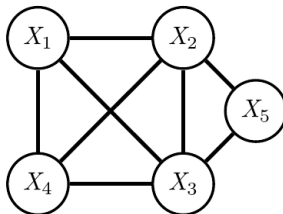
**Directed Acyclic Graph (DAG)** Graph in which by following the direction of the arrows a node will never be visited more than once.

**Parents and Children**  $X_i$  is a parent of  $X_j$  if there is a link from  $X_i$  to  $X_j$ .  $X_i$  is a child of  $X_j$  if there is a link from  $X_j$  to  $X_i$ .

**Ancestors and Descendants** The ancestors of a node  $X_i$  are the nodes with a directed path ending at  $X_i$ . The descendants of  $X_i$  are the nodes with a directed path beginning at  $X_i$ .

### 10.1.3 Undirected Graph

**Undirected Graph** Graph that all the edges are undirected.



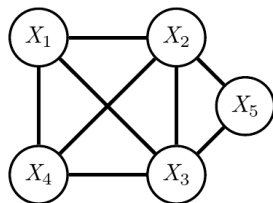
**Clique** A clique is a fully connected subset of nodes.  $(X_1, X_2, X_4)$  forms a (non-maximal) clique.

**Maximal Clique** Clique which is not a subset of a larger clique.  $(X_1, X_2, X_3, X_4)$  and  $(X_2, X_3, X_5)$  are both maximal cliques.

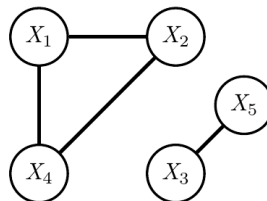
### 10.1.4 Connectivity

**Connected Graph** There is a path between every pair of vertices.

**Connected Components** In a non-connected graph, the connected components are the connected-subgraphs.  $(X_1, X_2, X_4)$  and  $(X_3, X_5)$  are the two connected components.



Connected Graph

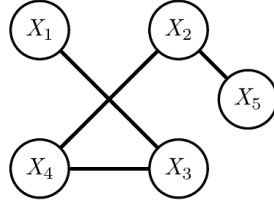


Connected Components

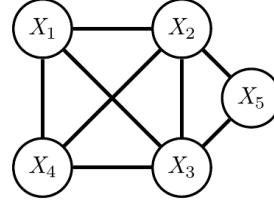
### 10.1.5 Connectedness

**Singly-connected** There is only one path from any node  $a$  to another node  $b$ .

**Multiply-connected** A graph is multiply-connected if it is not singly-connected.



Singly-connected

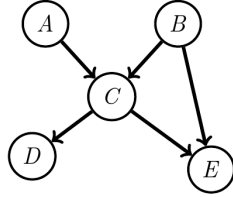


Multiply-connected

## 10.2 Belief Networks

### 10.2.1 Definition

A belief network is a directed acyclic graph in which each node is associated with the conditional probability of the node given its parents. The joint distribution is obtained by taking the product of the conditional probabilities.



$p(E|B, C)$

$$P(A, B, C, D, E) = P(A) P(B) P(C|A, B) P(D|C) P(E|B, C)$$

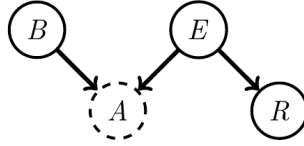
### 10.2.2 Uncertain Evidence

**Definition** In soft/uncertain evidence the variable is in more than one state, with the strength of our belief about each state being given by probabilities. For example, if  $y$  has the states  $\text{dom}(y) = \{\text{red}, \text{blue}, \text{green}\}$ , the vector  $(0.6, 0.1, 0.3)$  could represent the probabilities of the respective states.

**Hard Evidence** We are certain that a variable is in a particular state. In this state, all the probability mass is in one of the vector components,  $(0, 0, 1)$ .

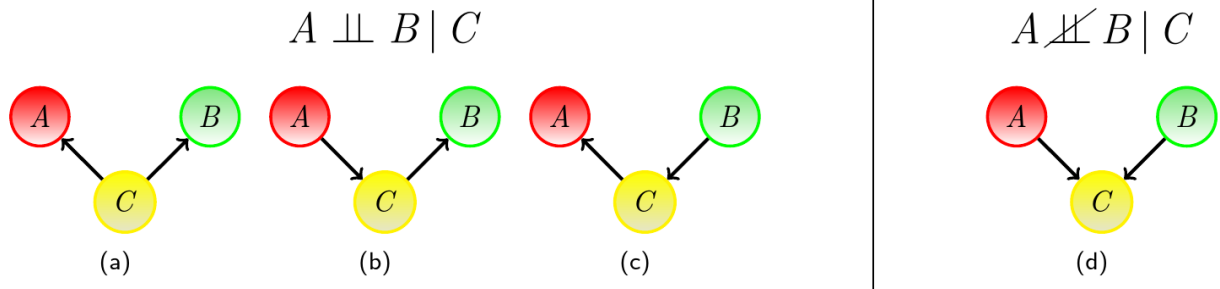
**Inference** Inference with soft-evidence can be achieved using Bayes' rule. Writing the soft-evidence as  $\tilde{y}$ , we have  $P(x|\tilde{y}) = \sum_y P(x|y) P(y|\tilde{y})$ , where  $P(y = i|\tilde{y})$  represents the probability that  $y$  is in state  $i$  under the soft-evidence.

**Jeffrey's Rule** For variables  $x, y$  and  $P_1(x, y)$ , how do we form a joint distribution given soft-evidence  $\tilde{y}$ ? (a) From the conditional we first define  $P_1(x|y) = \frac{P_1(x, y)}{\sum_x P_1(x, y)}$ . (b) Define the joint. The soft-evidence  $P(y|\tilde{y})$  then defines a new joint distribution  $P_2(x, y|\tilde{y}) = P_1(x|y)P_1(y|\tilde{y})$ . One can therefore view soft-evidence as defining a new joint distribution. We use a dashed circle to represent a variable in an uncertain state.



### 10.2.3 Independence

#### Conditionally Independent



In (a), (b) and (c), A, B are conditionally independent given C.

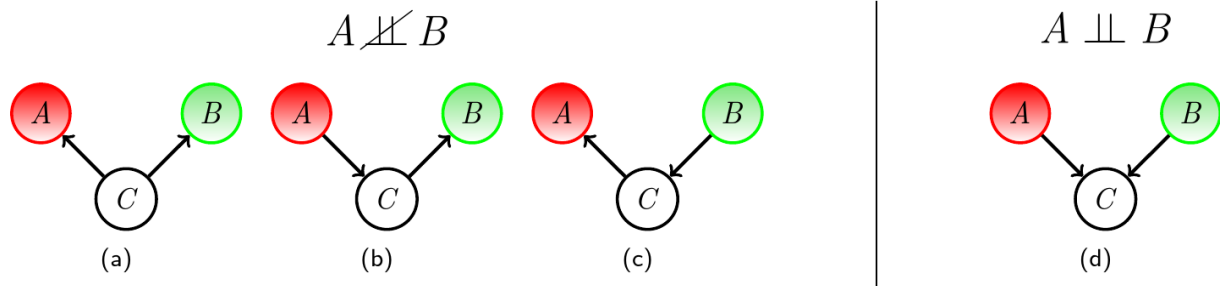
$$(a) \quad P(A, B|C) = \frac{P(A, B, C)}{P(C)} = \frac{P(A|C)P(B|C)P(C)}{P(C)} = P(A|C) P(B|C)$$

$$(b) \quad P(A, B|C) = \frac{P(A, B, C)}{P(C)} = \frac{P(A)P(C|A)P(B|C)}{P(C)} = \frac{P(A, C)P(B|C)}{P(C)} = P(A|C) P(B|C)$$

$$(c) \quad P(A, B|C) = \frac{P(A, B, C)}{P(C)} = \frac{P(A|C)P(C|B)P(B)}{P(C)} = \frac{P(A|C)P(B, C)}{P(C)} = P(A|C) P(B|C)$$

In (d) the variables A, B are conditionally dependent given C,  $P(A, B|C) \propto P(C|A, B) P(A) P(B)$ .

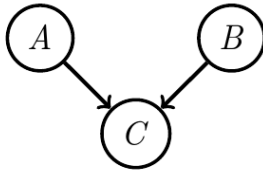
#### Marginally Dependent



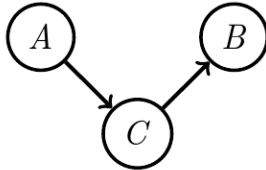
In (a), (b) and (c), the variables A, B are marginally dependent. In (d) the variables A, B are marginally independent.

$$P(A, B) = \sum_C P(A, B, C) = \sum_C P(A) P(B) P(C|A, B) = P(A) P(B)$$

### Colliders



If C has more than one incoming link, then  $A \not\perp\!\!\!\perp B|C$ . In this case C is called *collider*.



If C has at most one incoming link, then  $A \perp\!\!\!\perp B|C$  and  $A \not\perp\!\!\!\perp B$ . In this case C is called *non-collider*.

### 10.2.4 General Rule for Independence in Belief Networks

Given three sets of nodes  $\mathcal{X}$ ,  $\mathcal{Y}$ ,  $\mathcal{C}$ , if all paths from any element of  $\mathcal{X}$  to any element of  $\mathcal{Y}$  are blocked by  $\mathcal{C}$ , then  $\mathcal{X}$  and  $\mathcal{Y}$  are conditionally independent given  $\mathcal{C}$ . A path  $\mathcal{P}$  is blocked by  $\mathcal{C}$  if at least one of the following conditions is satisfied:

1. There is a collider in the path  $\mathcal{P}$  such that neither the collider nor any of its descendants is in the conditioning set  $\mathcal{C}$ .
2. There is a non-collider in the path  $\mathcal{P}$  that is in the conditioning set  $\mathcal{C}$ .

### Independence of $\mathcal{X}$ and $\mathcal{Y}$

When the conditioning set is empty  $\mathcal{C} = \emptyset$ , then a path  $\mathcal{P}$  from an element of  $\mathcal{X}$  to an element of  $\mathcal{Y}$  is blocked if there is a collider on the path. Hence  $\mathcal{X}$  and  $\mathcal{Y}$  are independent if every path from an element of  $\mathcal{X}$  to any element of  $\mathcal{Y}$  has a collider.

### d-connected

We use the term that  $\mathcal{X}$  and  $\mathcal{Y}$  are “d-connected” by  $\mathcal{Z}$  if there is any path from  $\mathcal{X}$  to  $\mathcal{Y}$  that is not blocked by  $\mathcal{Z}$ . If  $\mathcal{Z}$  is the empty set then we just say that  $\mathcal{X}$  and  $\mathcal{Y}$  are d-connected.

### Separation and Independence

Note first that d-separation and connection are properties of the graph (not of the distribution). d-separation implies that  $\mathcal{X} \perp\!\!\!\perp \mathcal{Y}|\mathcal{Z}$ , but d-connection does not necessarily imply conditional dependence. That is, for any distribution in which  $\mathcal{X}$  and  $\mathcal{Y}$  are “d-separated” by  $\mathcal{Z}$ , then no matter what the settings of the conditional tables are, then conditional independence holds, namely  $\mathcal{X} \perp\!\!\!\perp \mathcal{Y}|\mathcal{Z}$ .

## 10.2.5 Markov Equivalence

### Skeleton

Formed from a graph by removing the arrows.

### Immortality

An immortality in a DAG is a configuration of three nodes, A,B,C such that C is a child of both A and B, with A and B not directly connected.

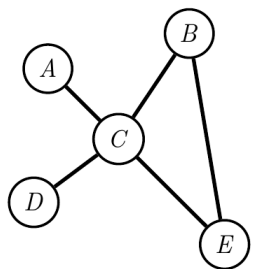
### Markov Equivalence

Markov equivalence Two graphs represent the same set of independence assumptions if and only if they have the same skeleton and the same set of immoralities.

## 10.3 Markov Networks

### 10.3.1 Definition

A Markov Network is an undirected graph in which there is a potential (non-negative function)  $\psi$  defined on each maximal clique.

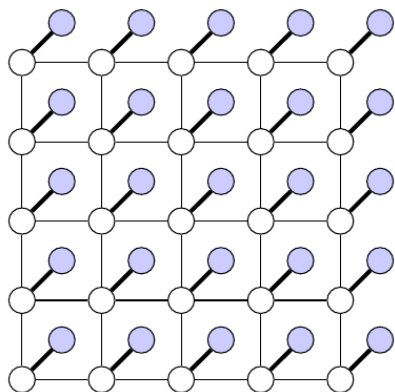


$$P(A, B, C, D, E) = \frac{1}{Z} \psi(A, C) \psi(C, D) \psi(B, C, E)$$

$$Z = \sum_{A, B, C, D, E} \psi(A, C) \psi(C, D) \psi(B, C, E)$$

### 10.3.2 Examples

#### Binary Image



$X = \{X_i, i = 1, \dots, D\}$   $X_i \in \{-1, 1\}$  : clean pixel

$Y = \{Y_i, i = 1, \dots, D\}$   $Y_i \in \{-1, 1\}$  : corrupted pixel

$\phi(Y_i, X_i) = e^{\gamma X_i Y_i}$  : encourage  $Y_i$  and  $X_i$  to be similar

$\psi(X_i, X_j) = e^{\beta X_i X_j}$  : encourage the image to be smooth

$$P(X, Y) \propto \left[ \prod_{i=1}^D \phi(Y_i, X_i) \right] \left[ \prod_{i \sim j} \psi(X_i, X_j) \right]$$

Boltzmann Machine

The Ising Model

**10.3.3 Independence**

**10.3.4 Expressiveness of Markov and Belief Networks**

**10.3.5 Factor Graphs**

**10.4 Markov Chains**

**10.5 Hidden Markov Models**

# Appendix A

## Statistical Assessment

### A.1 Hypothesis Testing

#### A.1.1 Testing Basics

**Null Hypothesis  $H_0$**  The hypothesis we would like to test.

**Alternative Hypothesis  $H_1$**  An alternative result when  $H_0$  is rejected. In most cases the alternative hypothesis is simply the negation of the null hypothesis.

**P-value** The p-value is the probability of observing a test statistic,  $X$ , as or more extreme than the value  $x$  seen in the data, under the assumption that the null hypothesis,  $H_0$ , is true. The p-value is most certainly not the probability of  $H_0$  being true.

**False Positives (Type I Error)** Rejecting  $H_0$  when it is true.

**False Negatives (Type II Error)** Not rejecting  $H_0$  when it is false. (N.B. Not rejecting  $H_0$  is not the same as accepting  $H_0$ )

**Power** The power of a hypothesis test is the probability of avoiding a false negative.

#### A.1.2 Testing Procedure

A common testing procedure includes the following steps:

1. Specify a null hypothesis ( $H_0$ ) and alternative hypothesis ( $H_1$ ).  
e.g.  
 $H_0 : \theta = 0.5$  The proportion of males and females is identical.  
 $H_1 : \theta < 0.5$  There is a smaller proportion of females than males.



2. Specify the level of the test.  
e.g. a common level=0.05
  - Bearing in mind the need to balance probabilities of Type I and Type II errors.
  - Reducing the level reduces the probability of a Type I error.
  - Increasing the level reduces the probability of a Type II error.
3. Specify a suitable test statistic.  
e.g.  $X = \text{The number of females} = 15$ .
4. Determine the distribution of the test statistic under  $H_0$ .  
e.g.  $X \sim \mathbf{Bin}(40, 0.5)$
5. Determine what it means to be “more extreme” by considering  $H_0$  and  $H_1$ .  
e.g.  $H_1 : \theta < 0.5$ , so smaller values of  $X$  are more extreme.
6. Determine the corresponding p-value.  
e.g.  $p = p(X \leq 15) = 0.077$
7. Reject  $H_0$  if the p-value is less than the level of the test.  
e.g.  $p > 0.05$ , so we fail to reject  $H_0$  in this instance. Conclude that the proportion of females and males is identical.

An alternative procedure is that rather than determining a p-value, we may determine a critical region for the test statistic – the set of all test statistic values which would cause us to reject  $H_0$ .

e.g. level = 0.05,  $p(X \leq 15) = 0.077$ ,  $p(X \leq 14) = 0.04 \implies CR = \{0, 1, 2, \dots, 14\}$ .

We may therefore simply compare our observed value to the critical region to judge whether to reject  $H_0$ .

### A.1.3 Power Investigation

### A.1.4 Useful Tests

## A.2 Confidence Intervals

## A.3 Bootstrap