

Mathematics, Statistics, Optimization, and Machine Learning

Sunghee Yun

December 28, 2019

Contents

I	Mathematics	7
1	Calculus	9
1.1	Basics	10
1.2	Multivariate functions	10
1.3	Chain rule	11
1.4	Integration	11
2	Convex analysis	13
2.1	Convex function	14
2.1.1	First order condition	14
2.1.2	Second order condition	17
3	Linear Algebra	19
3.1	Vector space	20
3.2	Eigenvalues	20
3.2.1	Basic definitions	20
3.2.2	Symmetric matrices	21
3.3	Positive definiteness	21
3.4	Matrix norms	22
II	Optimization	23
4	Convex Optimization	25
4.1	Mathematical optimization problem	26
4.2	Convex optimization problem	27
4.3	Duality	28
4.3.1	The Lagrange dual problem	28
4.3.1.1	Examples	28
4.3.2	Interpretations	29
4.3.2.1	Max-min characterization of weak and strong duality	29
4.3.2.2	Saddle-point interpretation	30
4.4	Unconstrained minimization	31
4.4.1	Gradient descent method	31

4.4.1.1	Examples	31
5	Portfolio optimization	33
5.1	Problem formulation	34
5.1.1	A portfolio optimization problem	34
III	Statistics	37
6	Statistics Basics	39
6.1	Correlation coefficients	39
6.2	Transformation of a random variable via a function	39
6.2.1	Scale random variable	39
6.2.2	Multivariate random variable	41
6.2.3	Data Examples	41
7	Various distributions	43
7.1	Log-normal distribution	43
7.1.1	Some statistics	44
7.1.2	Parameter estimation	46
8	Bayesian Statistics	47
8.1	Bayesian Theorem	47
8.2	Bayesian Inference	47
8.3	Conjugate prior	47
8.3.1	Bernoulli distribution	47
8.3.2	Gaussian distribution	48
9	Information Theory	49
9.1	Basics	49
9.1.1	Entropy	49
9.1.2	Mutual Information	49
9.1.3	Relative Entropy (Kullback–Leibler divergence)	49
IV	Machine Learning	51
10	Machine Learning Basics	53
10.1	Optimal Predictor	54
10.2	Bias and Variance	54
11	Optimization for Machine Learning	57
11.1	Gradient method	58
11.2	Stochastic gradient method	58
12	Bayesian Network	59

13 Collaborative Filtering	61
13.1 Item-based Collaborative Filtering	62
13.1.1 Rating matrix modeling for menu personalization for mobile shopping app	62
13.1.1.1 Rating matrix modeling	63
13.1.2 Value augmentation based on Bayesian MAP	64
13.1.3 Similarity measure among items	64
13.1.3.1 Cosine similarity	65
13.1.3.2 Cosine similarity when prior distribution is used	65
13.1.3.3 Correlation coefficient similarity	66
13.1.4 Data value transformation	66
13.1.4.1 TFIDF (or tf-idf)	67
13.1.4.2 Okapi BM25 transformation	67
13.1.5 Recommendation based on item similarities	67
13.2 Collaborative Filtering using Matrix Factorization	68
13.2.1 Problem definition and formulations	68
13.2.2 Solution methods	70
13.2.2.1 Matrix factorization via singular value decomposition (SVD)	70
13.2.2.2 Matrix factorization via gradient descent (GD) method	71
13.2.2.3 Matrix factorization via alternating gradient descent (GD) method	73
13.2.2.4 Matrix factorization via stochastic gradient descent (SGD) method	73
13.2.2.5 Matrix factorization via alternating least-squares (ALS)	74
13.2.2.6 Weighted matrix factorization via alternating least-squares (ALS)	76
13.2.3 Collaborative filtering for implicit feedback dataset	76
13.2.3.1 Regularization coefficient conversion	77
13.3 Appendix	79
13.3.1 Linear algebra	79
13.3.1.1 Singular value decomposition (SVD)	79
13.3.1.2 Singular value decomposition as rank- k approximation	80
14 Time Series Anomaly Detection	81
14.1 Real-Time Anomaly Detection	82
14.1.1 Computing Anomaly Likelihood	82
15 Reinforcement Learning	85
15.1 Finite Markov decision processes	86
15.1.1 Markov property	86
15.1.2 Policy	87
15.1.3 Return	87
15.1.4 State value function and action value function	88
15.2 Bellman equation	88
15.2.1 Bellman equations	88
15.2.2 Bellman optimality equations	89
15.3 Dynamic programming	90
15.3.1 Policy evaluation (prediction)	90
15.3.2 Policy iteration	91
15.3.3 Value iteration	91

15.4	Monte Carlo methods	91
15.4.1	Monte Carlo prediction	93
15.4.2	Monte Carlo control	94
15.4.3	Monte Carlo control without exploring starts	94
15.4.4	Off-policy prediction via important sampling	96
15.4.5	Off-policy Monte Carlo control	96
15.5	Temporal-difference learning	96
15.5.1	TD prediction	96
15.5.2	Sarsa: on-policy TD Control	98
15.5.3	Q-learning: off-policy TD control	98
15.5.4	Maximization bias and double learning	99
15.6	n -step bootstrapping	99
15.6.1	n -step TD prediction	100
15.6.2	n -step Sarsa	104
15.6.3	n -step off-policy learning	106
15.7	Planning and learning with tabular methods	106
15.7.1	Dyna: integrated planning, acting, and learning	106
15.8	On-policy Prediction with Approximation	106
15.9	On-policy Control with Approximation	106
15.10	Off-policy Methods with Approximation	109
15.11	Eligibility Traces	109
15.11.1	The λ -return	109
15.11.2	$TD(\lambda)$	110
15.11.3	Why $TD(\lambda)$ approximates the off-line λ -return algorithm?	111
15.11.4	Sarsa(λ)	116
15.11.5	Tabular methods using eligibility traces	116
15.12	Appendix: conditional probability and expected value	118

Part I

Mathematics

Chapter 1

Calculus

1.1 Basics

Theorem 1.1 (L'Hôpital's rule) Let $f : \mathbf{R} \rightarrow \mathbf{R}$ and $g : \mathbf{R} \rightarrow \mathbf{R}$ be differentiable on an open interval $I \subseteq \mathbf{R}$ except possibly at $c \in I$. If $\lim_{x \rightarrow c} f(x) = \lim_{x \rightarrow c} g(x) = 0$ or $\pm\infty$, $g'(x) \neq 0$ for all $x \in I \setminus \{c\}$, and $\lim_{x \rightarrow c} \frac{f'(x)}{g'(x)}$ exists, then

$$\lim_{x \rightarrow c} \frac{f(x)}{g(x)} = \lim_{x \rightarrow c} \frac{f'(x)}{g'(x)}. \quad (1.1)$$

Definition 1.1 (Taylor polynomial) Let $n \in \mathbf{Z}$ be a positive integer and let $f : \mathbf{R} \rightarrow \mathbf{R}$ be n times differentiable at $a \in \mathbf{R}$. The n -th order Taylor polynomial is defined by

$$\begin{aligned} T_{f,n}(x) &= f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \cdots + \frac{f^{(n)}(a)}{n!}(x-a)^n \\ &= \sum_{k=0}^n \frac{f^{(k)}(a)}{k!}(x-a)^k \end{aligned} \quad (1.2)$$

Theorem 1.2 (Taylor's theorem) Let $n \in \mathbf{Z}$ be a positive integer and let $f : \mathbf{R} \rightarrow \mathbf{R}$ be n times differentiable at $a \in \mathbf{R}$. Then there exists a function $h_n : \mathbf{R} \rightarrow \mathbf{R}$ such that

$$f(x) = T_{f,n}(x) + h_n(x)(x-a)^n \quad (1.3)$$

and $\lim_{x \rightarrow a} h_n(x) = 0$. The remainder is called the Peano form of the remainder.

Theorem 1.3 (Taylor's theorem) Let $n \in \mathbf{Z}$ be a positive integer, $a, b \in \mathbf{R}$, and $I_o = (a, b) \cup (b, a)$ and $I_c = [a, b] \cup [b, a]$. Let $f : \mathbf{R} \rightarrow \mathbf{R}$ be $n+1$ times differentiable on I_o and $f^{(n)}$ is continuous on I_c . Then for some $c \in I_o$,

$$f(b) = T_{f,n}(b) + \frac{f^{(n+1)}(c)}{(n+1)!}(b-a)^{n+1}. \quad (1.4)$$

The remainder is called the Peano form of the remainder.

1.2 Multivariate functions

Definition 1.2 (Jacobian matrix) Let $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$ be a differentiable function, i.e., the partial derivative $\partial f_j(x)/\partial x_i$ exists for every $1 \leq i \leq n$ and $1 \leq j \leq m$. Then the Jacobian matrix of f at x is defined by the function $D_f : \mathbf{R}^n \rightarrow \mathbf{R}^{m \times n}$ such that

$$D_f(x) = \begin{bmatrix} \partial f_1(x)/\partial x_1 & \partial f_1(x)/\partial x_2 & \cdots & \partial f_1(x)/\partial x_n \\ \partial f_2(x)/\partial x_1 & \partial f_2(x)/\partial x_2 & \cdots & \partial f_2(x)/\partial x_n \\ \vdots & \vdots & \ddots & \vdots \\ \partial f_m(x)/\partial x_1 & \partial f_m(x)/\partial x_2 & \cdots & \partial f_m(x)/\partial x_n \end{bmatrix} \in \mathbf{R}^{m \times n}. \quad (1.5)$$

1.3 Chain rule

Theorem 1.4 *Let $f : \mathbf{R} \rightarrow \mathbf{R}^n$ and $g : \mathbf{R}^n \rightarrow \mathbf{R}$ be differentiable. Then $h : \mathbf{R} \rightarrow \mathbf{R}$ such that $h(t) = g(f(t))$ is also differentiable and*

$$h'(t) = \sum_{i=1}^n f'_i(t) \frac{\partial g}{\partial x_i}(f(t)) = \nabla^T g(f(t))^T D_f(t)$$

for all $t \in \text{dom } f$.

Corollary 1.1 *Let $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$ and $g : \mathbf{R}^m \rightarrow \mathbf{R}^p$ be differentiable. Then define a function $h : \mathbf{R}^n \rightarrow \mathbf{R}^p$ such that $h(x) = g(f(x))$ for all $x \in \text{dom } f$. Then h is differentiable and*

$$Dh(x) = Dg(f(x))Df(x) \quad (1.6)$$

where $Df : \mathbf{R}^n \rightarrow \mathbf{R}^{m \times n}$, $Dg : \mathbf{R}^m \rightarrow \mathbf{R}^{p \times m}$, and $Dh : \mathbf{R}^n \rightarrow \mathbf{R}^{p \times n}$ are the Jacobian matrix functions of f , g , and h respectively.

Corollary 1.2 *Let $f : \mathbf{R}^n \rightarrow \mathbf{R}$ be differentiable. Then for some $A \in \mathbf{R}^{n \times m}$ and $b \in \mathbf{R}^n$, define $g : \mathbf{R}^m \rightarrow \mathbf{R}$ such that $g(y) = f(Ay + b)$. Then*

$$\nabla g(y) = A^T \nabla f(Ay + b). \quad (1.7)$$

Corollary 1.3 *Let $f : \mathbf{R}^n \rightarrow \mathbf{R}$ be twice differentiable. Then for some $A \in \mathbf{R}^{n \times m}$ and $b \in \mathbf{R}^n$, define $g : \mathbf{R}^m \rightarrow \mathbf{R}$ such that $g(y) = f(Ay + b)$. Then*

$$\nabla^2 g(y) = A^T \nabla^2 f(Ay + b)A. \quad (1.8)$$

1.4 Integration

Lemma 1.1 *Let $A \in \mathbf{R}^{n \times n}$ be a nonsingular matrix. Suppose that the following integral exists for some $C \subseteq \mathbf{R}^n$.*

$$\int_C f(x) dx \quad (1.9)$$

Chapter 2

Convex analysis

2.1 Convex function

A function $f : \mathbf{R}^n \rightarrow \mathbf{R}$ is a convex function if, for all $x, y \in \mathbf{dom} f$ and all $0 \leq \lambda \leq 1$,

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y). \quad (2.1)$$

Theorem 2.1 *Let $f : \mathbf{R}^n \rightarrow \mathbf{R}$. Then for some $x \in \mathbf{dom} f$ and $v \in \mathbf{R}^n$, define a function $g_{x,v}(t) : \mathbf{R} \rightarrow \mathbf{R}$ such that $g_{x,v}(t) = f(x + tv)$ with the domain, $\{t \in \mathbf{R} \mid x + tv \in \mathbf{dom} f\}$. Then f is a convex function iff $g_{x,v}$ is a convex function for any $x \in \mathbf{dom} f$ and any $v \in \mathbf{R}^n$.*

Proof: Suppose that f is a convex function. Then for any $x \in \mathbf{dom} f$ and $v \in \mathbf{R}^n$, for any $s, t \in \{t \in \mathbf{R} \mid x + tv \in \mathbf{dom} f\}$ and any $\lambda \in \mathbf{R}$ such that $0 \leq \lambda \leq 1$,

$$\begin{aligned} g_{x,v}(\lambda s + (1 - \lambda)t) &= f(x + (\lambda s + (1 - \lambda)t)v) \\ &= f(\lambda(x + sv) + (1 - \lambda)(x + tv)) \\ &\leq \lambda f(x + sv) + (1 - \lambda)f(x + tv) = \lambda g_{x,v}(s) + (1 - \lambda)g_{x,v}(t). \end{aligned}$$

Therefore $g_{x,v}$ is a convex function.

Now assume that $g_{x,v}(t) : \mathbf{R} \rightarrow \mathbf{R}$ is a convex function for any $x \in \mathbf{dom} f$ and $v \in \mathbf{R}^n$. Then for any $x, y \in \mathbf{dom} f$ and any $\lambda \in \mathbf{R}$ such that $0 \leq \lambda \leq 1$,

$$\begin{aligned} f((1 - \lambda)x + \lambda y) &= f(x + \lambda(y - x)) \\ &= g_{x, y-x}(\lambda) = g_{x, y-x}((1 - \lambda) \cdot 0 + \lambda \cdot 1) \leq (1 - \lambda)g_{x, y-x}(0) + \lambda g_{x, y-x}(1) \\ &= (1 - \lambda)f(x) + \lambda f(y), \end{aligned}$$

thus, f is a convex function.

2.1.1 First order condition

Theorem 2.2 *If a function $f : \mathbf{R} \rightarrow \mathbf{R}$ is differentiable, it is a convex function iff, for all $x, y \in \mathbf{dom} f$,*

$$f(y) \geq f(x) + f'(x)(y - x). \quad (2.2)$$

Proof: Suppose that f is a convex function. Then assume that $y > x$. Let $h \in \mathbf{R}$ be a positive number such that $h < y - x$. Then the definition of convexity implies that

$$f(x + h) \leq (1 - \lambda)f(x) + \lambda f(y)$$

where $\lambda = h/(y - x)$ since

$$(1 - \lambda)x + \lambda y = x + \lambda(y - x) = x + h.$$

Thus

$$f(x + h) - f(x) \leq \lambda(f(y) - f(x)) = \frac{h}{y - x}(f(y) - f(x)),$$

which implies

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h} \leq \frac{f(y) - f(x)}{y - x}.$$

Therefore

$$f(y) - f(x) \geq f'(x)(y - x),$$

hence (2.2) is true when $y > x$.

We can prove (2.2) is true when $y < x$ using the very same method. Assume that $x > y$. Let $h \in \mathbf{R}$ be a positive number such that $h < x - y$. Then the definition of convexity implies that

$$f(x - h) \leq (1 - \lambda)f(x) + \lambda f(y)$$

where $\lambda = h/(x - y)$ since

$$(1 - \lambda)x + \lambda y = x + \lambda(y - x) = x - h.$$

Thus

$$f(x) - f(x - h) \geq \lambda(f(x) - f(y)) = \frac{h}{x - y}(f(x) - f(y)),$$

which implies

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x) - f(x - h)}{h} \geq \frac{f(x) - f(y)}{x - y} = \frac{f(y) - f(x)}{y - x}.$$

Therefore

$$f(y) - f(x) \geq f'(x)(y - x),$$

hence (2.2) is true when $y < x$. It is obvious that (2.2) is true when $y = x$. Hence we have just proved that if $f : \mathbf{R} \rightarrow \mathbf{R}$ is a convex function, then (2.2) holds for any $x, y \in \mathbf{dom} f$.

Now we prove the converse. Suppose that (2.2) holds for any $x, y \in \mathbf{dom} f$. Now let $x, y \in \mathbf{dom} f$ and $\lambda \in \mathbf{R}$ such that $0 \leq \lambda \leq 1$. Let $z = \lambda x + (1 - \lambda)y$. Then (2.2) implies that

$$f(x) - f(z) \geq f'(z)(x - z) = (1 - \lambda)f'(z)(x - y) \quad (2.3)$$

and

$$f(y) - f(z) \geq f'(z)(y - z) = \lambda f'(z)(y - x) \quad (2.4)$$

If we multiply λ on both sides of (2.3), multiply $1 - \lambda$ on both sides of (2.4), and add both sides, we have

$$\lambda(f(x) - f(z)) + (1 - \lambda)(f(y) - f(z)) \geq \lambda f(x) + (1 - \lambda)f(y) - f(z) \geq 0,$$

hence

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y).$$

Therefore f is a convex function.

Corollary 2.1 *Let $f : \mathbf{R} \rightarrow \mathbf{R}$ be differentiable. Then f is a convex function iff the derivative of f is a nondecreasing function.*

Proof: Suppose that f is a convex function. Let $x, y \in \mathbf{dom} f$ such that $x < y$. Then Theorem 2.2 implies

$$f(y) \geq f(x) + f'(x)(y - x)$$

and

$$f(x) \geq f(y) + f'(y)(x - y),$$

thus

$$f'(x) \leq \frac{f(y) - f(x)}{y - x} = \frac{f(x) - f(y)}{x - y} \leq f'(y)$$

since $y > x$. Therefore f' is a nondecreasing function.

Now we prove the converse. Suppose that f' is a nondecreasing function. Then if $y > x$, the mean value theorem implies that there exists some $z \in (x, y)$ such that

$$\frac{f(y) - f(x)}{y - x} = f'(z).$$

Since f' is nondecreasing, we have

$$f'(x) \leq \frac{f(y) - f(x)}{y - x} \leq f'(y),$$

thus

$$f(y) \geq f(x) + f'(x)(y - x) \tag{2.5}$$

and

$$f(y) \leq f(x) + f'(y)(y - x). \tag{2.6}$$

Therefore (2.5) implies that f satisfies (2.2). Now if $x > y$, (2.6) implies that

$$f(x) \leq f(y) + f'(x)(x - y) \Leftrightarrow f(y) \geq f(x) + f'(x)(y - x)$$

which again implies that f satisfies (2.2). Therefore (2.2) implies that f is a convex function.

Corollary 2.2 *If a function $f : \mathbf{R}^n \rightarrow \mathbf{R}$ is differentiable, it is a convex function iff, for all $x, y \in \mathbf{dom} f$,*

$$f(y) \geq f(x) + \nabla f(x)^T(y - x). \tag{2.7}$$

Proof: Suppose that f is a convex function. Let $x, y \in \mathbf{dom} f$. If we let $g_{x,y-x} : \mathbf{R} \rightarrow \mathbf{R}$ be a function such that $g_{x,y-x}(t) = f(x + t(y - x))$, Theorem 2.1 implies $g_{x,y-x}$ is a convex function. Therefore Theorem 2.2 together with Corollary 1.2 implies

$$f(y) = g_{x,y-x}(1) \geq g_{x,y-x}(0) + g'_{x,y-x}(0)(1 - 0) = f(x) + \nabla f(x)^T(y - x)$$

for any $x, y \in \mathbf{dom} f$.

Now suppose that (2.7) holds for any $x, y \in \mathbf{dom} f$. Then Corollary 1.2 implies that, for any $r, s \in \mathbf{R}$ and $v \in \mathbf{R}^n$ such that $r, s \in \{t \in \mathbf{R} \mid x + tv \in \mathbf{dom} f\}$,

$$g_{x,v}(r) = f(x + rv) \geq f(x + sv) + (r - s)\nabla f(x + sv)^T v = g_{x,v}(s) + g'_{x,v}(s)(r - s).$$

Thus Theorem 2.2 implies $g_{x,v}$ is a convex function for any $x \in \mathbf{dom} f$ and $v \in \mathbf{R}^n$. Therefore by Theorem 2.1, f is a convex function.

2.1.2 Second order condition

Theorem 2.3 *If a function $f : \mathbf{R} \rightarrow \mathbf{R}$ is twice differentiable, it is a convex function iff, for all $x \in \mathbf{dom} f$,*

$$f''(x) \geq 0. \quad (2.8)$$

Proof: Suppose that f is a convex function. Then Corollary 2.1 implies that f' is a nondecreasing function, hence

$$f''(x) = \lim_{h \rightarrow 0} \frac{f'(x+h) - f'(x)}{h} = \lim_{h \rightarrow 0^+} \frac{f'(x+h) - f'(x)}{h} \geq 0.$$

Now if $f''(x) \geq 0$ for all $x \in \mathbf{dom} f$, the mean value theorem implies that f' is a nondecreasing function.

Theorem 2.4 *If a function $f : \mathbf{R}^n \rightarrow \mathbf{R}$ is twice differentiable, it is a convex function iff, for all $x \in \mathbf{dom} f$,*

$$\nabla^2 f(x) \succeq 0. \quad (2.9)$$

Proof: Suppose that f is a convex function. Then Theorem 2.1 implies that for any $x \in \mathbf{dom} f$ and any $v \in \mathbf{R}^n$, the function $g_{x,v} : \mathbf{R} \rightarrow \mathbf{R}$ such that $g_{x,v}(t) = f(x + tv)$ is a convex function in $\{t \in \mathbf{R} \mid x + tv \in \mathbf{dom} f\}$. Then Theorem 2.3 together with Corollary 1.3 implies that

$$v^T \nabla^2 f(x) v = g''_{x,v}(0) \geq 0.$$

Therefore $\nabla^2 f(x) \succeq 0$ for any $x \in \mathbf{dom} f$.

Now if $\nabla^2 f(x) \succeq 0$ for all $x \in \mathbf{dom} f$, then Corollary 1.3 implies that $g''_{x,v}(t) = v^T \nabla^2 f(x + tv) v \geq 0$ for any $x \in \mathbf{dom} f$ and any $v \in \mathbf{R}^n$. Then Theorem 2.3 implies $g_{x,v}$ is a convex function for any $x \in \mathbf{dom} f$ and any $v \in \mathbf{R}^n$, hence by Theorem 2.1, f is a convex function.

Chapter 3

Linear Algebra

3.1 Vector space

Cauchy–Schwarz inequality: for $a, b \in \mathbf{C}^n$,

$$|a^H b| \leq \|a\|_2 \|b\|_2. \quad (3.1)$$

The generalized form: for $f : [0, 1] \rightarrow \mathbf{C}$ and $g : [0, 1] \rightarrow \mathbf{C}$,

$$\left| \int_0^1 \overline{f(t)} g(t) dt \right| \leq \left(\int_0^1 |f(t)|^2 dt \right)^{1/2} \left(\int_0^1 |g(t)|^2 dt \right)^{1/2}. \quad (3.2)$$

Hölder's inequality: for $a, b \in \mathbf{C}^n$, $p > 1$, and $q > 1$ such that $1/p + 1/q = 1$,

$$|a^H b| \leq \|a\|_p \|b\|_q. \quad (3.3)$$

The generalized form: for $f : [0, 1] \rightarrow \mathbf{C}$, $g : [0, 1] \rightarrow \mathbf{C}$, $p > 1$, and $q > 1$ such that $1/p + 1/q = 1$,

$$\left| \int_0^1 \overline{f(t)} g(t) dt \right| \leq \left(\int_0^1 |f(t)|^p dt \right)^{1/p} \left(\int_0^1 |g(t)|^q dt \right)^{1/q}. \quad (3.4)$$

When $b = \mathbf{1}$, the Cauchy–Schwarz inequality implies

$$\|a\|_1 \leq n \|a\|_2 \quad (3.5)$$

3.2 Eigenvalues

3.2.1 Basic definitions

Given a square matrix $A \in \mathbf{R}^{n \times n}$, if there exist $\lambda \in \mathbf{C}$ and nonzero $v \in \mathbf{C}^n$ such that

$$Av = \lambda v \quad (3.6)$$

then λ is called an eigenvalue of A and v is called an eigenvector associated with λ .

If there exist n linearly independent eigenvectors, we have

$$A \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix} = \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix} \mathbf{diag}(\lambda_1, \dots, \lambda_n) \quad (3.7)$$

or

$$AV = V\Lambda \quad (3.8)$$

where

$$V = \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix} \in \mathbf{C}^{n \times n} \quad (3.9)$$

and

$$\Lambda = \mathbf{diag}(\lambda_1, \dots, \lambda_n) = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix} \in \mathbf{C}^{n \times n}. \quad (3.10)$$

In this case, A is said to be diagonalizable.

Since V is nonsingular, *i.e.*, invertible, we can rewrite (3.8) as

$$A = V\Lambda V^{-1} \Leftrightarrow V^{-1}AV = \Lambda. \quad (3.11)$$

3.2.2 Symmetric matrices

Given a symmetric matrix $A = A^T \in \mathbf{R}^{n \times n}$, all the eigenvalues are real and we can choose n real orthonormal eigenvectors, *i.e.*, we can find n eigenvectors $v_1, \dots, v_n \in \mathbf{R}^n$ associated with n eigenvalues, $\lambda_1, \dots, \lambda_n \in \mathbf{R}$ such that

$$\|v_i\| = 1 \quad (3.12)$$

for $i = 1, \dots, n$ and

$$v_i^T v_j = 0 \quad (3.13)$$

for $1 \leq i \neq j \leq n$. Thus, all symmetric matrices are diagonalizable.

Now (3.11) becomes

$$A = V\Lambda V^T \Leftrightarrow V^T A V = \Lambda \quad (3.14)$$

since

$$V^T V = I_n \quad (3.15)$$

where $I_n \in \mathbf{R}^{n \times n}$ is the identity matrix. We can rewrite (3.14) as

$$A = \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix} \mathbf{diag}(\lambda_1, \dots, \lambda_n) \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix} = \sum_{i=1}^n \lambda_i v_i v_i^T. \quad (3.16)$$

3.3 Positive definiteness

- A symmetric matrix $A = A^T \in \mathbf{R}^{n \times n}$ is called positive semidefinite if for all $x \in \mathbf{R}^n$,

$$x^T A x \geq 0. \quad (3.17)$$

- A symmetric matrix $A = A^T \in \mathbf{R}^{n \times n}$ is called positive definite if for all nonzero $x \in \mathbf{R}^n$,

$$x^T A x > 0. \quad (3.18)$$

- The set of all the n -by- n positive semidefinite matrices is (sometimes) denoted by \mathcal{S}_+^n , *i.e.*,

$$\mathcal{S}_+^n = \{A = A^T \in \mathbf{R}^{n \times n} \mid x^T A x \geq 0 \text{ for all } x \in \mathbf{R}^n\}. \quad (3.19)$$

- The set of all the n -by- n positive definite matrices is (sometimes) denoted by \mathcal{S}_{++}^n , *i.e.*,

$$\mathcal{S}_{++}^n = \{A = A^T \in \mathbf{R}^{n \times n} \mid x^T A x > 0 \text{ for all nonzero } x \in \mathbf{R}^n\}. \quad (3.20)$$

- $A = A^T \in \mathbf{R}^{n \times n}$ is positive semidefinite if and only if all the eigenvalues of A are nonnegative.

- $A = A^T \in \mathbf{R}^{n \times n}$ is positive definite if and only if all the eigenvalues of A are positive.

Proof: For symmetric $A = A^T$, there exist orthogonal $V \in \mathbf{R}^{n \times n}$ and diagonal $\Lambda \in \mathbf{R}^{n \times n}$ such that

$$A = V\Lambda V^T = \sum_{i=1}^n \lambda_i v_i v_i^T,$$

thus for any $x \in \mathbf{R}^n$,

$$x^T A x = x^T \left(\sum_{i=1}^n \lambda_i v_i v_i^T \right) x = \sum_{i=1}^n \lambda_i x^T v_i v_i^T x = \sum_{i=1}^n \lambda_i (v_i^T x)^2.$$

Therefore if all λ_i are nonnegative, $x^T A x \geq 0$ for any $x \in \mathbf{R}^n$, hence $A \in \mathcal{S}_+^n$. Now assume $A \in \mathcal{S}_+^n$, but $\lambda_j < 0$ for some $j \in \{1, \dots, n\}$. Then

$$v_j^T A v_j = \sum_{i=1}^n \lambda_i (v_i^T v_j)^2 = \sum_{i=1}^n \lambda_i \delta_{i,j} = \lambda_j < 0 \quad (3.21)$$

since v_1, \dots, v_n are orthonormal where $\delta_{i,j}$ is the [Kronecker delta function](#), hence $A \notin \mathcal{S}_+^n$. Therefore if $A \in \mathcal{S}_+^n$, all λ_i are nonnegative.

Therefore $A \in \mathcal{S}_+^n$ if and only if all λ_i are nonnegative.

Now assume that all λ_i are positive. Then for all nonzero $x \in \mathbf{R}^n$, there exists $i \in \{1, \dots, n\}$ such that $v_i^T x \neq 0$ since if $v_i^T x = 0$ for all i , then $V^T x = 0$, hence $x = 0$ since V is nonsingular. Therefore

$$x^T A x = \sum_{i=1}^n \lambda_i (v_i^T x)^2 \geq \lambda_j (v_j^T x)^2 > 0. \quad (3.22)$$

Thus, $A \in \mathcal{S}_{++}^n$.

Now assume that $A \in \mathcal{S}_{++}^n$. If $\lambda_j \leq 0$ for some $j \in \{1, \dots, n\}$, then

$$v_j^T A v_j = \sum_{i=1}^n \lambda_i \delta_{i,j} = \lambda_j \leq 0, \quad (3.23)$$

hence $A \notin \mathcal{S}_{++}^n$. Therefore if $A \in \mathcal{S}_{++}^n$, all λ_i are positive.

Therefore $A \in \mathcal{S}_{++}^n$ if and only if all λ_i are positive.

3.4 Matrix norms

$$\text{dist}(C_{\text{org1}}, C_{\text{org2}}) = \|C_{\text{org1}} - C_{\text{org2}}\| = |\lambda_{\max}(C_{\text{org1}} - C_{\text{org2}})| \quad (3.24)$$

Part II

Optimization

Chapter 4

Convex Optimization

4.1 Mathematical optimization problem

A mathematical optimization problem can be expressed as

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0 \text{ for } i = 1, \dots, m \\ & && h_i(x) = 0 \text{ for } i = 1, \dots, p \end{aligned} \quad (4.1)$$

where $x \in \mathbf{R}^n$ is the optimization variable, $f_0 : \mathbf{R}^n \rightarrow \mathbf{R}$ is the objective function, $f_i : \mathbf{R}^n \rightarrow \mathbf{R}$ for $i = 1, \dots, m$ are the inequality constraint functions, and $h_i : \mathbf{R}^n \rightarrow \mathbf{R}$ for $i = 1, \dots, p$ are the equality constraint functions.

The conditions, $f_i(x) \leq 0$ for $i = 1, \dots, m$, are called inequality constraints and the conditions, $h_i(x) = 0$ for $i = 1, \dots, p$ are called equation constraints.

Note that this formulation covers pretty much every single-objective optimization problem. For example, consider the following optimization problem.

$$\begin{aligned} & \text{maximize} && f(x_1, x_2) \\ & \text{subject to} && x_1 \geq x_2 \\ & && x_1 + x_2 = 2 \end{aligned} \quad (4.2)$$

This problem can be cast into an equivalent problem as follows.

$$\begin{aligned} & \text{minimize} && -f(x_1, x_2) \\ & \text{subject to} && -x_1 + x_2 \leq 0 \\ & && x_1 + x_2 - 2 = 0 \end{aligned} \quad (4.3)$$

The feasible set for (4.1) is defined by the set of $x \in \mathbf{R}^n$ which satisfies all the constraints. Also, the optimal value for (4.1) is the infimum of $f_0(x)$ while x is in the feasible set. When the infimum is achievable, we define the optimal solution set as the set of all feasible x achieving the infimum value. These are defined in mathematically rigorous terms below.

- The feasible set for (4.1) is defined by

$$\mathcal{F} = \{x \in \mathcal{D} \mid f_i(x) \leq 0 \text{ for } i = 0, \dots, m, h_j(x) = 0 \text{ for } j = 1, \dots, p\} \subseteq \mathbf{R}^n \quad (4.4)$$

where

$$\mathcal{D} = \left(\bigcap_{0 \leq i \leq m} \text{dom } f_i \right) \cap \left(\bigcap_{1 \leq i \leq p} \text{dom } h_i \right). \quad (4.5)$$

- The optimal value for (4.1) is defined by

$$p^* = \inf_{x \in \mathcal{F}} f_0(x) \quad (4.6)$$

We use the conventions that $p^* = -\infty$ if $f_0(x)$ is unbounded below for $x \in \mathcal{F}$ and that $p^* = \infty$ if $\mathcal{F} = \emptyset$.

- The optimal solution set for (4.1) is defined by

$$\mathcal{X}^* = \{x \in \mathcal{F} \mid f_0(x) = p^*\}. \quad (4.7)$$

4.2 Convex optimization problem

A mathematical optimization problem is called a convex optimization problem if the objective function and all the inequality constraint functions are convex functions and all the equality constraint functions are affine functions.

Hence, a convex optimization problem can be expressed as

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0 \text{ for } i = 1, \dots, m \\ & && Ax = b \end{aligned} \tag{4.8}$$

where $x \in \mathbf{R}^n$ is the optimization variable, $f_i : \mathbf{R}^n \rightarrow \mathbf{R}$ for $i = 0, \dots, n$ are convex functions, $h_i : \mathbf{R}^n \rightarrow \mathbf{R}$ for $i = 1, \dots, p$ are the equality constraint functions, $A \in \mathbf{R}^{p \times n}$, and $b \in \mathbf{R}^p$.

A function, $f : \mathbf{R}^n \rightarrow \mathbf{R}$, is called a convex function if $\text{dom } f \subseteq \mathbf{R}^n$ is a convex set and for all $x, y \in \text{dom } f$ and all $0 \leq \lambda \leq 1$,

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \tag{4.9}$$

where $\text{dom } f \subseteq \mathbf{R}^n$ denotes the domain of f .

A convex optimization enjoys a number of nice theoretical and practical properties.

- A local minimum of a convex optimization problem is a global minimum, *i.e.*, if for some $R > 0$ and $x_0 \in \mathcal{F}$, $\|x - x_0\| < R$ and $x \in \mathcal{F}$ imply $f_0(x_0) \leq f_0(x)$, then $f_0(x_0) \leq f_0(x)$ for all $x \in \mathcal{F}$.

Proof: Assume that $x_0 \in \mathcal{F}$ is a local minimum, *i.e.*, for some $R > 0$, $\|x - x_0\| < R$ and $x \in \mathcal{F}$ imply $f_0(x_0) \leq f_0(x)$.

Now assume that x_0 is not a global minimum, *i.e.*, there exists $y \in \mathcal{F}$ such that $y \neq x_0$ and $f_0(y) < f_0(x_0)$. Then for $z = \lambda y + (1 - \lambda)x_0$ with $\lambda = \min\{R/\|y - x_0\|, 1\}/2$, the convexity of f_0 implies

$$f_0(z) \leq \lambda f_0(y) + (1 - \lambda)f_0(x_0) \tag{4.10}$$

since $0 < \lambda \leq 1/2 < 1$. Furthermore

$$\|z - x_0\| = \lambda\|y - x_0\| \leq R/2, \tag{4.11}$$

hence $f_0(z) \geq f_0(x_0)$, which together with (4.10) implies

$$f_0(x_0) \leq f_0(z) \leq \lambda f_0(y) + (1 - \lambda)f_0(x_0) < \lambda f_0(x_0) + (1 - \lambda)f_0(x_0) = f_0(x_0), \tag{4.12}$$

which is a contradiction. Therefore there is no $y \in \mathcal{F}$ such that $y \neq x_0$ and $f_0(y) < f_0(x_0)$. Therefore x_0 is a global minimum.

- For a unconstrained problem, *i.e.*, the problem (4.8) with $m = p = 0$, with differentiable objective function, $x \in \text{dom } f_0$ is an optimal solution if and only if $\nabla f_0(x) = 0 \in \mathbf{R}^n$.

Proof: The Taylor theorem implies that for any $x, y \in \mathbf{dom} f_0$,

$$f_0(y) = f_0(x) + \nabla f_0(x)^T(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f_0(z)(y - x) \quad (4.13)$$

for some z on the line segment having x and y as its end points, *i.e.*, $z = \alpha x + (1 - \alpha)y$ for some $0 \leq \alpha \leq 1$. Since $\nabla^2 f_0(x) \succeq 0$ for any $z \in \mathbf{dom} f_0$, we have

$$f_0(y) \geq f_0(x) + \nabla f_0(x)^T(y - x) \quad (4.14)$$

Thus, if for some $x_0 \in \mathbf{R}^n$, $\nabla f_0(x_0) = 0$, for any $x \in \mathbf{dom} f_0$,

$$f_0(x) \geq f_0(x_0) + \nabla f_0(x_0)^T(x - x_0) = f_0(x_0), \quad (4.15)$$

hence x_0 is an optimal solution. Now assume that x_0 is an optimal solution, but $\nabla f_0(x_0) \neq 0$. Then for any $k > 0$, if we let $x = x_0$ and $y = x_0 - k \nabla f_0(x_0)$, (4.13) becomes

$$\begin{aligned} f_0(y) &= f_0(x_0) + \nabla f_0(x_0)^T(-k \nabla f_0(x_0)) + \frac{k^2}{2} \nabla f_0(x_0)^T \nabla^2 f_0(z) \nabla f_0(x_0) \\ &= f_0(x_0) - k \|\nabla f_0(x_0)\|^2 + \frac{k^2}{2} \nabla f_0(x_0)^T \nabla^2 f_0(z) \nabla f_0(x_0) \end{aligned}$$

for all $y = x_0 - k \nabla f_0(x_0) \in \mathbf{dom} f_0$.

Since for $k < 2\|\nabla f_0(x_0)\|^2 / \nabla f_0(x_0)^T \nabla^2 f_0(z) \nabla f_0(x_0)$, $-k\|\nabla f_0(x_0)\|^2 + \frac{k^2}{2} \nabla f_0(x_0)^T \nabla^2 f_0(z) \nabla f_0(x_0) < 0$, thus $f_0(y) < f_0(x_0)$, hence the contradiction. Therefore, if x_0 is an optimal solution for the unconstrained problem, $\nabla f_0(x_0) = 0$.

4.3 Duality

4.3.1 The Lagrange dual problem

4.3.1.1 Examples

4.3.1.1.1 Standard form LP

$$\begin{aligned} &\text{minimize} && c^T x \\ &\text{subject to} && Ax = b \\ &&& x \succeq 0 \end{aligned} \quad (4.16)$$

The Lagrange dual problem is

$$\begin{aligned} &\text{maximize} && -b^T \nu \\ &\text{subject to} && A^T \nu + c \geq 0 \end{aligned} \quad (4.17)$$

4.3.1.1.2 Inequality form LP

$$\begin{aligned} & \text{minimize} && c^T x \\ & \text{subject to} && Ax \preceq b \end{aligned} \tag{4.18}$$

The Lagrange dual problem is

$$\begin{aligned} & \text{maximize} && -b^T \lambda \\ & \text{subject to} && A^T \lambda + c = 0 \\ & && \lambda \succeq 0 \end{aligned} \tag{4.19}$$

4.3.1.1.3 Least-squares solution of linear equations

$$\begin{aligned} & \text{minimize} && (1/2)x^T x \\ & \text{subject to} && Ax = b \end{aligned} \tag{4.20}$$

The Lagrange dual problem is

$$\text{maximize} \quad -(1/2)\nu^T A A^T \nu - b^T \nu \tag{4.21}$$

4.3.1.1.4 Entropy maximization

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^n x_i \log x_i \\ & \text{subject to} && Ax = b \\ & && \mathbf{1}^T x = 1 \end{aligned} \tag{4.22}$$

with domain $\mathcal{D} = \mathbf{R}_+^n$

The Lagrange dual problem is

$$\begin{aligned} & \text{maximize} && -b^T \lambda - \log \left(\sum_{i=1}^n \exp(-a_i^T \lambda) \right) \\ & \text{subject to} && \lambda \succeq 0 \end{aligned} \tag{4.23}$$

4.3.2 Interpretations**4.3.2.1 Max-min characterization of weak and strong duality**

We first note that for any $f : X \times Y \rightarrow \mathbf{R}$, we have

$$\sup_{y \in Y} \inf_{x \in X} f(x, y) \leq \inf_{x \in X} \sup_{y \in Y} f(x, y). \tag{4.24}$$

This inequality is called *max-min inequality*.

We can prove this as follows. Let $g : Y \rightarrow \mathbf{R}$ be a function defined by $g(y) = \inf_{x \in X} f(x, y)$ and let $h : X \rightarrow \mathbf{R}$ be a function defined by $h(x) = \sup_{y \in Y} f(x, y)$. Then we have that for any $x \in X$ and $y \in Y$

$$g(y) = \inf_{x \in X} f(x, y) \leq f(x, y), \tag{4.25}$$

which implies that for any $x \in X$

$$\sup_{y \in Y} g(y) \leq \sup_{y \in Y} f(x, y) = h(x). \quad (4.26)$$

This again implies that

$$\sup_{y \in Y} g(y) \leq \inf_{x \in X} h(x), \quad (4.27)$$

hence the proof.

4.3.2.2 Saddle-point interpretation

Suppose $f : X \times Y \rightarrow \mathbf{R}$. We refer a point $(\tilde{x}, \tilde{y}) \in X \times Y$ a *saddle-point* for f (and X and Y) if

$$f(\tilde{x}, y) \leq f(\tilde{x}, \tilde{y}) \leq f(x, \tilde{y}) \quad (4.28)$$

for all $x \in X$ and $y \in Y$.

Now if x^* and λ^* are primal and dual optimal points for a problem in which strong duality obtains, the form a saddle-point for the Lagrangian. Conversely, if (x, λ) is a saddle-point of the Lagrangian, then x is primal optimal, λ is dual optimal, and the optimal duality gap is zero.

To prove these, assume that $x^* \in \mathcal{D}$ and $(\lambda^*, \nu^*) \in \mathbf{R}_+^m \times \mathbf{R}^p$ are primal and dual optimal points for a problem in which strong duality obtains. Then for any $x \in \mathcal{D}$ and $(\lambda, \nu) \in \mathbf{R}_+^m \times \mathbf{R}^p$, we have

$$L(x^*, \lambda, \nu) = f_0(x^*) + \sum_{i=1}^m \lambda_i f_i(x^*) + \sum_{i=1}^p \nu_i h_i(x^*) \leq f_0(x^*) = g(\lambda^*, \nu^*) \leq L(x, \lambda^*, \nu^*) \quad (4.29)$$

where the left inequality comes from the fact that $\lambda_i f_i(x^*) \leq 0$ for all $i = 1, \dots, m$ and $h_i(x^*) = 0$ for all $i = 1, \dots, p$ and the right inequality comes from the definition of (Lagrange) dual function. Now from the complementary slackness we know that $\lambda_i f_i(x^*) = 0$ for all $i = 1, \dots, m$. Therefore

$$L(x^*, \lambda^*, \nu^*) = f_0(x^*), \quad (4.30)$$

thus we have

$$L(x^*, \lambda, \nu) \leq L(x^*, \lambda^*, \nu^*) \leq L(x, \lambda^*, \nu^*), \quad (4.31)$$

hence the proof.

Now suppose that $\tilde{x} \in \mathcal{D}$ and $(\tilde{\lambda}, \tilde{\nu}) \in \mathbf{R}_+^m \times \mathbf{R}^p$ are the saddle-point of the Lagrangian, *i.e.*, for all $x \in \mathcal{D}$ and $(\lambda, \nu) \in \mathbf{R}_+^m \times \mathbf{R}^p$,

$$L(\tilde{x}, \lambda, \nu) \leq L(\tilde{x}, \tilde{\lambda}, \tilde{\nu}) \leq L(x, \tilde{\lambda}, \tilde{\nu}). \quad (4.32)$$

First we show that \tilde{x} is a feasible point. The left inequality says that for all $(\lambda, \nu) \in \mathbf{R}_+^m \times \mathbf{R}^p$,

$$L(\tilde{x}, \lambda, \nu) = f_0(\tilde{x}) + \sum_{i=1}^m \lambda_i f_i(\tilde{x}) + \sum_{i=1}^p \nu_i h_i(\tilde{x}) \leq L(\tilde{x}, \tilde{\lambda}, \tilde{\nu}) \quad (4.33)$$

If $f_i(\tilde{x}) > 0$ for some $i \in \{1, \dots, m\}$ or $h_i(\tilde{x}) \neq 0$ for some $i \in \{1, \dots, p\}$, $L(\tilde{x}, \lambda, \nu)$ is unbounded above and the above inequality cannot hold. Therefore $f_i(\tilde{x}) \leq 0$ for all $i \in \{1, \dots, m\}$ and $h_i(\tilde{x}) = 0$

for all $i \in \{1, \dots, p\}$, *i.e.*, \tilde{x} is primal feasible. Since the inequality must hold when $\lambda = 0$ and $\nu = 0$, we have

$$f(\tilde{x}) \leq L(\tilde{x}, \tilde{\lambda}, \tilde{\nu}). \quad (4.34)$$

The right inequality of (4.32) implies that

$$L(\tilde{x}, \tilde{\lambda}, \tilde{\nu}) \leq g(\tilde{\lambda}, \tilde{\nu}) = \inf_{x \in \mathcal{D}} L(x, \tilde{\lambda}, \tilde{\nu}), \quad (4.35)$$

which implies that $f_0(\tilde{x}) \leq g(\tilde{\lambda}, \tilde{\nu})$. Since $g(\lambda, \nu)$ is an underestimator of $f_0(x)$ for any feasible $x \in \mathcal{D}$ and $(\tilde{\lambda}, \tilde{\nu}) \in \mathbf{R}_+^m \times \mathbf{R}^p$, *i.e.*, $g(\tilde{\lambda}, \tilde{\nu}) \leq f_0(\tilde{x})$, thus $g(\tilde{\lambda}, \tilde{\nu}) = f_0(\tilde{x})$. Therefore \tilde{x} is an optimal solution for the primal problem and $(\tilde{\lambda}, \tilde{\nu})$ is an optimal solution for the dual problem, hence the proof.

4.4 Unconstrained minimization

4.4.1 Gradient descent method

4.4.1.1 Examples

4.4.1.1.1 A quadratic problem in \mathbf{R}^2 We consider the quadratic objective function on \mathbf{R}^2

$$f(x) = \frac{1}{2}(x_1^2 + \gamma x_2^2) \quad (4.36)$$

where $\gamma > 0$.

We apply the gradient descent method with exact line search. The gradient of f is

$$\nabla f(x) = \begin{bmatrix} x_1 \\ \gamma x_2 \end{bmatrix} \quad (4.37)$$

Let $\tilde{f} : \mathbf{R}_+ \rightarrow \mathbf{R}$ defined by $\tilde{f}(t) = f(x - t\nabla f(x))$. Now

$$\tilde{f}(t) = f\left(\begin{bmatrix} (1-t)x_1 \\ (1-\gamma t)x_2 \end{bmatrix}\right) = \frac{1}{2}((1-t)^2 x_1^2 + \gamma(1-\gamma t)^2 x_2^2) \quad (4.38)$$

and

$$\frac{d}{dt}\tilde{f}(t) = -(1-t)x_1^2 - \gamma^2(1-\gamma t)x_2^2 = 0 \quad (4.39)$$

implies

$$t = \frac{x_1^2 + \gamma^2 x_2^2}{x_1^2 + \gamma^3 x_2^2} \quad (4.40)$$

minimizes $\tilde{f}(t)$. Since

$$1-t = \frac{\gamma^2(\gamma-1)x_2^2}{x_1^2 + \gamma^3 x_2^2} \quad (4.41)$$

and

$$1-\gamma t = \frac{(1-\gamma)x_1^2}{x_1^2 + \gamma^3 x_2^2} \quad (4.42)$$

Thus the exact line search yields

$$x^+ = x - t\nabla f(x) = \begin{bmatrix} (1-t)x_1 \\ (1-\gamma t)x_2 \end{bmatrix} = \frac{(1-\gamma)x_1x_2}{x_1^2 + \gamma^3x_2^2} \begin{bmatrix} -\gamma^2x_2 \\ x_1 \end{bmatrix}. \quad (4.43)$$

If $x = \alpha[\gamma \ 1]^T$, then

$$x^+ = \frac{\alpha^3(1-\gamma)\gamma}{\alpha^2\gamma^2(1+\gamma)} \begin{bmatrix} -\gamma^2 \\ \gamma \end{bmatrix} = \alpha \frac{1-\gamma}{1+\gamma} \begin{bmatrix} -\gamma \\ 1 \end{bmatrix}. \quad (4.44)$$

If $x = \alpha[-\gamma \ 1]^T$, then

$$x^+ = -\frac{\alpha^3(1-\gamma)\gamma}{\alpha^2\gamma^2(1+\gamma)} \begin{bmatrix} -\gamma^2 \\ -\gamma \end{bmatrix} = \alpha \frac{1-\gamma}{1+\gamma} \begin{bmatrix} \gamma \\ 1 \end{bmatrix}. \quad (4.45)$$

Therefore if $x^{(0)} = [\gamma \ 1]^T$, then

$$x^{(k)} = \left(\frac{1-\gamma}{1+\gamma}\right)^k \begin{bmatrix} (-1)^k\gamma \\ 1 \end{bmatrix} = \left(\frac{\gamma-1}{\gamma+1}\right)^k \begin{bmatrix} \gamma \\ (-1)^k \end{bmatrix}. \quad (4.46)$$

Chapter 5

Portfolio optimization

5.1 Problem formulation

Suppose that we have n assets to invest on and that the return of each asset per unit invest is modeled by random variables R_i for $i = 1, \dots, n$. Then we want to decide the amount of investment on each asset, $x_i \in \mathbf{R}$ for $i = 1, \dots, n$, so that it optimizes the overall investment (in certain senses).

For formulation, we use the following definitions.

- Define a vector random variable $R \in \mathbf{R}^n$ such that

$$R = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} \in \mathbf{R}^n. \quad (5.1)$$

- Let $r \in \mathbf{R}^n$ be the expected value of R , *i.e.*,

$$r = \mathbf{E}(R) = \begin{bmatrix} \mathbf{E}(R_1) \\ \vdots \\ \mathbf{E}(R_n) \end{bmatrix} = \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix} \in \mathbf{R}^n. \quad (5.2)$$

- Define a vector $x \in \mathbf{R}$ which is an aggregate of the investments:

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbf{R}^n. \quad (5.3)$$

- Define a feasible set $\mathcal{X} \subseteq \mathbf{R}^n$ for x . For example, if we have a limit on the total investment,

$$\mathcal{X} = \{x \in \mathbf{R}^n \mid \sum_{i=1}^n c_i x_i \leq c_{\max}\}, \quad (5.4)$$

or if we have the minimum and maximum amount to invest for each asset, we'd have

$$\mathcal{X} = \{x \in \mathbf{R}^n \mid d_{\min} \leq x_i \leq d_{\max} \text{ for } i = 1, \dots, n\}. \quad (5.5)$$

Generally, we'd prefer \mathcal{X} to be a convex set, *i.e.*, for any $x, y \in \mathcal{X}$ and $0 \leq \lambda \leq 1$,

$$\lambda x + (1 - \lambda)y \in \mathcal{X}. \quad (5.6)$$

5.1.1 A portfolio optimization problem

A portfolio optimization problem can be formulized by

$$\begin{aligned} &\text{maximize} && f(x) = \mathbf{E}(Z) \\ &\text{minimize} && g(x) = \mathbf{Var}(Z) \\ &\text{subject to} && x \in \mathcal{X} \end{aligned} \quad (5.7)$$

where the optimization variable is $x \in \mathbf{R}^n$ and

$$Z = \sum_{i=1}^n x_i R_i = x^T R \quad (5.8)$$

where $\mathbf{E}(\cdot)$ and $\mathbf{Var}(\cdot)$ refer to the expected value and the variance operators respectively.

This problem formulation tries to *maximize the expected return* while *minimizing the variance or uncertainty or risk*, which generally makes sense.

(5.4) (5.5)

Note that

$$\mathbf{E}(Z) = \mathbf{E}(x^T R) = \mathbf{E}\left(\sum_{i=1}^n x_i R_i\right) = \sum_{i=1}^n x_i \mathbf{E}(R_i) = \sum_{i=1}^n x_i r_i = r^T x \quad (5.9)$$

and

$$\begin{aligned} \mathbf{Var}(Z) &= \mathbf{E}(Z - \mathbf{E}(Z))^2 = \mathbf{E}(x^T R - x^T r)^2 \\ &= \mathbf{E}(x^T (R - r))^2 = \mathbf{E}(x^T (R - r)(R - r)^T x) \\ &= x^T \mathbf{E}(R - r)(R - r)^T x = x^T \Sigma_R x \end{aligned}$$

where $\Sigma_R = \mathbf{E}(R - r)(R - r)^T$ is the [covariance matrix](#) of R . Note that $\Sigma_R \in \mathcal{S}_+^n$ since for any $y \in \mathbf{R}^n$,

$$y^T \Sigma_R y = y^T \mathbf{E}(R - r)(R - r)^T y = \mathbf{E}(y^T (R - r))^2 \geq 0. \quad (5.10)$$

Thus, (5.7) can be rewritten as

$$\begin{aligned} &\text{maximize} && f(x) = r^T x \\ &\text{minimize} && g(x) = x^T \Sigma_R x \\ &\text{subject to} && x \in \mathcal{X}. \end{aligned} \quad (5.11)$$

Part III

Statistics

Chapter 6

Statistics Basics

6.1 Correlation coefficients

The correlation coefficients of two random variables, X and Y , is defined by

$$\rho_{X,Y} = \frac{\mathbf{E}(X - \mu_X)(Y - \mu_Y)}{\sqrt{\mathbf{E}(X - \mu_X)^2 \mathbf{E}(Y - \mu_Y)^2}} \quad (6.1)$$

6.2 Transformation of a random variable via a function

6.2.1 Scale random variable

Assume two random variables, $X \in \mathbf{R}$ and $Y \in \mathbf{R}$, which are related by a function $g : \mathbf{R} \rightarrow \mathbf{R} \in C^1$ such that

$$Y = g(X). \quad (6.2)$$

Now let's derive an equation for the probability density function (PDF) of Y given the PDF of X , $f_X : \mathbf{R} \rightarrow \mathbf{R}_+$.

The definition of cumulative distribution function (CDF) of Y implies that

$$F_Y(y) = \mathbf{Prob}\{Y \leq y\} = \mathbf{Prob}\{g(X) \leq y\} \quad (6.3)$$

for any $y \in \mathbf{R}$.

Now if we assume that g is a strictly increasing function, it has its inverse function $g^{-1} : g(\mathbf{R}) \rightarrow \mathbf{R}$ and (6.3) becomes

$$F_Y(y) = \mathbf{Prob}\{g(X) \leq y\} = \mathbf{Prob}\{X \leq g^{-1}(y)\} = F_X(g^{-1}(y))$$

for any $y \in g(\mathbf{R})$. Thus, we can differentiate both sides to have

$$f_Y(y) = \frac{d}{dy} F_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y) = \frac{1}{g'(g^{-1}(y))} f_X(g^{-1}(y)), \quad (6.4)$$

since (1.4) implies that

$$1 = \frac{d}{dx}g(g^{-1}(x)) = g'(g^{-1}(x))\frac{d}{dx}g^{-1}(x), \quad (6.5)$$

i.e., the derivative of the inverse function is the inverse of the derivative of the original function.

Now if we assume that g is a strictly decreasing function, we have

$$F_Y(y) = \mathbf{Prob}\{g(X) \leq y\} = \mathbf{Prob}\{x \geq g^{-1}(y)\} = 1 - F_X(g^{-1}(y)) + \mathbf{Prob}\{x = g^{-1}(y)\},$$

and

$$f_Y(y) = \frac{d}{dy}F_Y(y) = -\frac{d}{dy}F_X(g^{-1}(y)) = -f_X(g^{-1}(y))\frac{d}{dy}g^{-1}(y) = -\frac{1}{g'(g^{-1}(y))}f_X(g^{-1}(y)). \quad (6.6)$$

Since $g'(y) > 0$ for a strictly increasing function, and $g'(y) < 0$ for a strictly decreasing function, (6.4) and (6.6) imply

$$f_Y(y) = \frac{1}{|g'(g^{-1}(y))|}f_X(g^{-1}(y)) \quad (6.7)$$

for both cases.

Now consider a general function, g , *i.e.*, not necessarily a monotonic function. Suppose that $y \in g(\mathbf{R})$. Then for every $x \in \mathbf{R}$ such that $f(x) = y$, if $f'(x) \neq 0$, then $f(x)$ is locally strictly monotonic, *i.e.*, there exists $\delta > 0$ such that $f(x)$ is strictly monotonic for $x \in (x - \delta, x + \delta)$, hence (6.7) holds for such x .

The probability around y is a summation of the probabilities around such points, *i.e.*, the probabilities around all x such that $f(x) = y$ and $f'(x) \neq 0$. Therefore, we have

$$f_Y(y) = \sum_{x:g(x)=y} \frac{1}{|g'(x)|}f_X(x). \quad (6.8)$$

There is another way to derive the same equation (in a less strict way) which helps get more insight. Let's again suppose that g is a strictly increasing function. Then consider the probability that X lies in $(x, x + \Delta x)$. The probability should be the same as Y lies in $(y, y + \Delta y)$ where $y = g(x)$ and $\Delta y = g(x + \Delta x) - g(x)$, *i.e.*,

$$\begin{aligned} f_X(x)\Delta x &\approx \int_x^{x+\Delta x} f_X(x) dx = \mathbf{Prob}\{x \leq X \leq x + \Delta x\} \\ &= \mathbf{Prob}\{y \leq Y \leq y + \Delta y\} = \int_y^{y+\Delta y} f_Y(y) dy \approx f_Y(y)\Delta y. \end{aligned}$$

The approximation becomes the equality when Δx goes to 0. Therefore we have

$$f_Y(y) = \lim_{\Delta x \rightarrow 0} \frac{\Delta x}{\Delta y} f_X(x) = \frac{1}{\lim_{\Delta x \rightarrow 0} \frac{g(x+\Delta x) - g(x)}{\Delta x}} f_X(x) = \frac{1}{g'(x)} f_X(x), \quad (6.9)$$

which is equivalent to (6.4). Following the very same argument as before will lead to (6.8), *i.e.*, applying the same method to strictly decreasing case, *etc.*

6.2.2 Multivariate random variable

6.2.3 Data Examples

Suppose that we have n random variables, X_1, \dots, X_n and they are independent and identically distributed Gaussian, $\mathcal{N}(0, 1)$. Then assume that a random variable, Y , is the sum of the X_i 's, *i.e.*,

$$Y = \sum_{i=1}^n X_i = X_1 + \dots + X_n \quad (6.10)$$

Then the covariance of X_i and Y for each i is

$$\mathbf{Cov}(X_i, Y) = \mathbf{E}(X_i - \mathbf{E}X_i)(Y - \mathbf{E}Y) = \mathbf{E} \left(\sum_{j=1}^n X_i X_j \right) = 1 \quad (6.11)$$

and the variance of Y is

$$\mathbf{Var}(Y) = \mathbf{E}(Y - \mathbf{E}Y)^2 = \mathbf{E} \left(\sum_{j=1}^n X_i X_j \right)^2 = \sum_{i=1}^n \mathbf{E}X_i^2 = n. \quad (6.12)$$

Hence, the correlation coefficient of X_i and Y for each i is

$$\rho_{X_i, Y} = \mathbf{Cov}(X_i, Y) / \sqrt{\mathbf{Var}X_i \mathbf{Var}Y} = 1/\sqrt{n}. \quad (6.13)$$

Therefore Y has clear relation with X_i 's, but each correlation coefficient can be arbitrarily small as n grows!

Chapter 7

Various distributions

7.1 Log-normal distribution

We say Y is log-normally distributed, if, for $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$,

$$Y = \exp(X). \tag{7.1}$$

Then (??) implies that

$$\begin{aligned} f_Y(y) &= \frac{1}{\exp(\log(y))} \cdot \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(\log(y) - \mu_X)^2}{2\sigma_X^2}\right) \\ &= \frac{1}{\sqrt{2\pi}\sigma_X y} \exp\left(-\frac{(\log(y) - \mu_X)^2}{2\sigma_X^2}\right). \end{aligned} \tag{7.2}$$

7.1.1 Some statistics

The definition of the expected value implies

$$\begin{aligned}
\mathbf{E}Y &= \int_0^\infty y f_Y(y) dy = \int_0^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(\log(y) - \mu_X)^2}{2\sigma_X^2}\right) dy \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - \mu_X)^2}{2\sigma_X^2}\right) \exp(x) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - \mu_X)^2}{2\sigma_X^2} + x\right) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{x^2 - 2(\mu_X + \sigma_X^2)x + \mu_X^2}{2\sigma_X^2}\right) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - (\mu_X + \sigma_X^2))^2 + \mu_X^2 - (\mu_X + \sigma_X^2)^2}{2\sigma_X^2}\right) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - (\mu_X + \sigma_X^2))^2 - 2\mu_X\sigma_X^2 - \sigma_X^4}{2\sigma_X^2}\right) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - (\mu_X + \sigma_X^2))^2 - \sigma_X^2(2\mu_X + \sigma_X^2)}{2\sigma_X^2}\right) dx \\
&= \exp\left(\frac{2\mu_X + \sigma_X^2}{2}\right) \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - (\mu_X + \sigma_X^2))^2}{2\sigma_X^2}\right) dx \\
&= \exp\left(\frac{2\mu_X + \sigma_X^2}{2}\right)
\end{aligned}$$

since $dy = \exp(x)dx$ and $\frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - (\mu_X + \sigma_X^2))^2}{2\sigma_X^2}\right)$ is the PDF of a random variable $\sim \mathcal{N}(\mu_X + \sigma_X^2, \sigma_X^2)$, thus

$$\mu_Y = \mathbf{E}Y = \exp\left(\frac{2\mu_X + \sigma_X^2}{2}\right). \quad (7.3)$$

Similarly,

$$\begin{aligned}
\mathbf{E}Y^2 &= \int_0^\infty y^2 f_Y(y) dy = \int_0^\infty \frac{1}{\sqrt{2\pi}\sigma_X} y \exp\left(-\frac{(\log(y) - \mu_X)^2}{2\sigma_X^2}\right) dy \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp(x) \exp\left(-\frac{(x - \mu_X)^2}{2\sigma_X^2}\right) \exp(x) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - \mu_X)^2}{2\sigma_X^2} + 2x\right) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{x^2 - 2(\mu_X + 2\sigma_X^2)x + \mu_X^2}{2\sigma_X^2}\right) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - (\mu_X + 2\sigma_X^2))^2 + \mu_X^2 - (\mu_X + 2\sigma_X^2)^2}{2\sigma_X^2}\right) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - (\mu_X + 2\sigma_X^2))^2 - 4\mu_X\sigma_X^2 - 4\sigma_X^4}{2\sigma_X^2}\right) dx \\
&= \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - (\mu_X + 2\sigma_X^2))^2 - 4\sigma_X^2(\mu_X + \sigma_X^2)}{2\sigma_X^2}\right) dx \\
&= \exp(2(\mu_X + \sigma_X^2)) \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x - (\mu_X + 2\sigma_X^2))^2}{2\sigma_X^2}\right) dx \\
&= \exp(2(\mu_X + \sigma_X^2)),
\end{aligned}$$

thus

$$\sigma_Y^2 = \mathbf{Var}(Y) = \mathbf{E}Y^2 - (\mathbf{E}Y)^2 = \exp(2(\mu_X + \sigma_X^2)) - \exp(2\mu_X + \sigma_X^2) = (\exp(\sigma_X^2) - 1) \exp(2\mu_X + \sigma_X^2). \quad (7.4)$$

Note that (7.3) implies that

$$\mu_Y^2 = \exp(2\mu_X + \sigma_X^2), \quad (7.5)$$

hence

$$\sigma_Y^2 = (\exp(\sigma_X^2) - 1) \mu_Y^2. \quad (7.6)$$

This multiplicative dependency of the standard deviation on the expected value is attributed to the fact that $\log(Y) \sim \mathcal{N}(\mu_X, \sigma_X^2)$, *i.e.*, the log-scale of Y follows the normal distribution.

Now if we differentiate the PDF with respect to y , (7.2) implies that

$$\begin{aligned}
\frac{d}{dy} f_Y(y) &= \frac{d}{dy} \left(\frac{1}{\sqrt{2\pi}\sigma_X y} \exp\left(-\frac{(\log(y) - \mu_X)^2}{2\sigma_X^2}\right) \right) \\
&= -\frac{1}{\sqrt{2\pi}\sigma_X y^2} \exp\left(-\frac{(\log(y) - \mu_X)^2}{2\sigma_X^2}\right) \\
&\quad + \frac{1}{\sqrt{2\pi}\sigma_X y} \exp\left(-\frac{(\log(y) - \mu_X)^2}{2\sigma_X^2}\right) \left(-\frac{(\log(y) - \mu_X)}{\sigma_X^2}\right) \frac{1}{y} \\
&= -\frac{1}{\sqrt{2\pi}\sigma_X y^2} \exp\left(-\frac{(\log(y) - \mu_X)^2}{2\sigma_X^2}\right) \left(1 + \frac{(\log(y) - \mu_X)}{\sigma_X^2}\right) \\
&= -\frac{1}{\sqrt{2\pi}\sigma_X^3 y^2} \exp\left(-\frac{(\log(y) - \mu_X)^2}{2\sigma_X^2}\right) (\log(y) - (\mu_X - \sigma_X^2)).
\end{aligned}$$

Equating the derivative to zero yields

$$y = \exp(\mu_X - \sigma_X^2), \quad (7.7)$$

which is the mode of Y .

7.1.2 Parameter estimation

Now assume that we have a log-normally distributed random variable, $Y \in \mathbf{R}_{++}$, with μ_Y and σ_Y^2 as its mean and variance. We derived the parameters of the source distribution, μ_X and σ_X .

The two equations, (7.3) and (7.4), imply

$$\begin{aligned} \mu_Y &= \exp(\mu_X + \sigma_X^2/2), \\ \sigma_Y^2 &= (\exp(\sigma_X^2) - 1) \exp(2\mu_X + \sigma_X^2) = (\exp(\sigma_X^2) - 1) \mu_Y^2, \end{aligned}$$

thus

$$\begin{aligned} \sigma_X^2 &= \log(1 + \sigma_Y^2/\mu_Y^2), \\ \mu_X &= \log(\mu_Y) - \sigma_X^2/2 = \log(\mu_Y) - \log(1 + \sigma_Y^2/\mu_Y^2)/2 = \frac{1}{2} \log \left(\frac{\mu_Y^2}{1 + \sigma_Y^2/\mu_Y^2} \right). \end{aligned}$$

Chapter 8

Bayesian Statistics

8.1 Bayesian Theorem

Suppose that A and B are two events with $P(B) \neq 0$. Then

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}. \quad (8.1)$$

This is called *Bayesian theorem*.

8.2 Bayesian Inference

$$p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)} = \frac{p(X|\theta)p(\theta)}{\int p(X|\theta)p(\theta)d\theta} \propto p(X|\theta)p(\theta) = \text{likelihood} \times \text{prior} \quad (8.2)$$

8.3 Conjugate prior

8.3.1 Bernoulli distribution

For Bernoulli distribution, the conjugate prior is the beta distribution.

$$p(\theta) = \frac{1}{B(a, b)} \theta^{a-1} (1 - \theta)^{b-1} \quad (8.3)$$

Then the posterior probability can be expressed as

$$\begin{aligned} p(\theta|X) &\propto p(X|\theta)p(\theta) \propto \left(\prod_{i=1}^n \theta^{I(x_i=1)} (1 - \theta)^{I(x_i=0)} \right) \theta^{a-1} (1 - \theta)^{b-1} \\ &= \theta^{k+a-1} (1 - \theta)^{n-k+b-1} \end{aligned}$$

where k is the number of 1s in X . Thus $p(\theta|X) \sim \text{Beta}(k+a, n-k+b)$, *i.e.*,

$$p(\theta|X) = \frac{1}{B(k+a, n-k+b)} \theta^{k+a-1} (1-\theta)^{n-k+b-1}. \quad (8.4)$$

8.3.2 Gaussian distribution

Suppose that $X \sim \mathcal{N}(\mu, \tau^{-1})$ and $\mu \sim \mathcal{N}(m, \lambda^{-1})$. Then the posterior distribution is

$$p(\mu|X) \sim \mathcal{N}\left(\frac{\tau \sum_{i=1}^n x_i + \lambda m}{n\tau + \lambda}, (n\tau + \lambda)^{-1}\right) \quad (8.5)$$

Chapter 9

Information Theory

9.1 Basics

9.1.1 Entropy

9.1.2 Mutual Information

The mutual information (MI) is defined by

$$I(X; Y) = \mathbf{E} \log \frac{f_{X,Y}(X, Y)}{f_X(X)f_Y(Y)} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \log \frac{f_{X,Y}(x, y)}{f_X(x)f_Y(y)} dx dy \quad (9.1)$$

9.1.3 Relative Entropy (Kullback–Leibler divergence)

The relative entropy of the two random distributions, p and q , is defined by

$$D(p\|q) = \mathbf{E}_p \log \left(\frac{p(X)}{q(X)} \right). \quad (9.2)$$

The relative entropy represents the difference measure between the two distributions. Unlike the mutual information, relative entropy is asymmetric, *i.e.*, in general, $D(p\|q) \neq D(q\|p)$. It is always nonnegative quantity since the Jensen's inequality implies that

$$-D(p\|q) = \mathbf{E}_p \log \left(\frac{q(X)}{p(X)} \right) \leq \log \left(\mathbf{E}_p \left(\frac{q(X)}{p(X)} \right) \right) = \log \left(\int_{-\infty}^{\infty} q(x) dx \right) = 0. \quad (9.3)$$

The relative entropy can be rewritten as

$$D(p\|q) = -\mathbf{E}_p \log q(X) + \mathbf{E}_p \log p(X) = -\mathbf{E}_p \log q(X) - H(p) \quad (9.4)$$

where the first term, $-\mathbf{E}_p \log q(X)$ is called the *cross-entropy* of p and q .

Part IV

Machine Learning

Chapter 10

Machine Learning Basics

10.1 Optimal Predictor

Consider a regression problem where we predict $Y \in \mathbf{R}^m$ given $X \in \mathbf{R}^n$. We want to design a predictor $g : \mathbf{R}^n \rightarrow \mathbf{R}^m$ so that $g(X) \sim Y$ in some statistical sense. We first show that $g(X) = \mathbf{E}(Y|X)$ is the optimal predictor (or estimator) in least-mean-square sense.

We define $g^* : \mathbf{R}^n \rightarrow \mathbf{R}^m$ where $g^*(x) = \mathbf{E}(Y|X = x)$. Then

$$\begin{aligned}
 \mathbf{E}\|g(X) - Y\|_2^2 &= \mathbf{E}\|g(X) - g^*(X) + g^*(X) - Y\|_2^2 \\
 &= \mathbf{E}\|g(X) - g^*(X)\|_2^2 + \mathbf{E}\|g^*(X) - Y\|_2^2 + 2\mathbf{E}(g(X) - g^*(X))^T(g^*(X) - Y) \\
 &= \mathbf{E}\|g(X) - g^*(X)\|_2^2 + \mathbf{E}\|g^*(X) - Y\|_2^2 + 2\mathbf{E}_X \mathbf{E}_Y ((g(X) - g^*(X))^T(g^*(X) - Y)|X) \\
 &= \mathbf{E}\|g(X) - g^*(X)\|_2^2 + \mathbf{E}\|g^*(X) - Y\|_2^2 + 2\mathbf{E}_X (g(X) - g^*(X))^T \mathbf{E}_Y (g^*(X) - Y|X) \\
 &= \mathbf{E}\|g(X) - g^*(X)\|_2^2 + \mathbf{E}\|g^*(X) - Y\|_2^2 + 2\mathbf{E}_X (g(X) - g^*(X))^T (g^*(X) - \mathbf{E}(Y|X)) \\
 &= \mathbf{E}\|g(X) - g^*(X)\|_2^2 + \mathbf{E}\|g^*(X) - Y\|_2^2 \geq \mathbf{E}\|g^*(X) - Y\|_2^2.
 \end{aligned}$$

Therefore $g^*(X)$ is the optimal predictor for Y in the least-mean-square sense.

10.2 Bias and Variance

In §10.1, we proved that $g^*(X) = \mathbf{E}(Y|X)$ is the optimal predictor (or estimator) in the least-mean-square sense. However, unless we have the full knowledge of the joint probability distribution of X and Y , *i.e.*, $p(X, Y)$, or know $\mathbf{E}(Y|X = x)$ as a function of x , it is not possible to obtain g^* .

Here we assume that we obtain the predictor for Y given X from a dataset \mathcal{D} where

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\} \subseteq \mathbf{R}^n \times \mathbf{R}^m. \quad (10.1)$$

Now suppose that we have a predictor $g(\cdot; \mathcal{D}) : \mathbf{R}^n \rightarrow \mathbf{R}^m$, which depends on \mathcal{D} .

¹Note that strictly speaking, \mathcal{D} is *not* a set since the order of $(x_i, y_i) \in \mathbf{R}^n \times \mathbf{R}^m$ matters, *i.e.*, if the order is changed, we generally have different predictor, and we are allowed to have identical data point. Thus, we should say \mathcal{D} is a (ordered) list of points, $(x_i, y_i) \in \mathbf{R}^n \times \mathbf{R}^m$.

Then the mean square error of this predictor can be decomposed as following.

$$\begin{aligned}
\mathbf{E}_{X,Y,\mathcal{D}} \|g(X; \mathcal{D}) - Y\|_2^2 &= \mathbf{E}_{X,Y,\mathcal{D}} \|g(X; \mathcal{D}) - g^*(X) + g^*(X) - Y\|_2^2 \\
&= \mathbf{E}_{X,Y,\mathcal{D}} \|g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y,\mathcal{D}} \|g^*(X) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,Y,\mathcal{D}} (g(X; \mathcal{D}) - g^*(X))^T (g^*(X) - Y) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,\mathcal{D}} \mathbf{E}_Y ((g(X; \mathcal{D}) - g^*(X))^T (g^*(X) - Y) | X, \mathcal{D}) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,\mathcal{D}} (g(X; \mathcal{D}) - g^*(X))^T \mathbf{E}_Y (g^*(X) - Y | X, \mathcal{D}) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,\mathcal{D}} (g(X; \mathcal{D}) - g^*(X))^T \mathbf{E}_Y (g^*(X) - Y | X) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,\mathcal{D}} (g(X; \mathcal{D}) - g^*(X))^T (g^*(X) - \mathbf{E}_Y(Y | X)) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) + \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2 \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D})\|_2^2 + \mathbf{E}_{X,\mathcal{D}} \|\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,\mathcal{D}} (g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}))^T (\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D})\|_2^2 + \mathbf{E}_X \|\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_X \mathbf{E}_{\mathcal{D}} ((g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}))^T (\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)) | X) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D})\|_2^2 + \mathbf{E}_X \|\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_X \mathbf{E}_{\mathcal{D}} (g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) | X)^T (\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D})\|_2^2 + \mathbf{E}_X \|\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_X (\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) | X)^T (\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D})\|_2^2 + \mathbf{E}_X \|\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)\|_2^2 + \mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2.
\end{aligned}$$

Note that we use the fact that $\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D})$ is a function of X only (hence does not depend on X).

In the last equation, the first term is called the *variance* since it is the expected value (with respect to X) of variance of the predictor, $g(X; \mathcal{D})$, with respect to the dataset, \mathcal{D} . It represents the extent to which the prediction varies around its expected value. The second term is the expected value of the square of the bias where the bias is defined to be the difference between the expected value of prediction with respect to dataset and the optimal prediction. The second term itself is sometimes called *bias*. The third term is called *noise* since it is caused by the intrinsic noise residing in Y which cannot be reduced even with the optimal predictor (in least-mean-square sense).

The following equation summarizes these three quantities.

$$\begin{aligned}
&\mathbf{E}_{X,Y,\mathcal{D}} \|g(X; \mathcal{D}) - Y\|_2^2 \\
&= \underbrace{\mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}} g(X; \mathcal{D})\|_2^2}_{\text{variance}} + \underbrace{\mathbf{E}_X \|\mathbf{E}_{\mathcal{D}} g(X; \mathcal{D}) - g^*(X)\|_2^2}_{\text{bias}} + \underbrace{\mathbf{E}_{X,Y} \|g^*(X) - Y\|_2^2}_{\text{noise}}.
\end{aligned} \tag{10.2}$$

In general, we do not know the optimal predictor; if we knew, we would not need to train our in

the first place. Thus we can only estimate $g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D})$ and $\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y$. The mean square error can also be expressed in these two quantities as follows.

$$\begin{aligned}
\mathbf{E}_{X,Y,\mathcal{D}} \|g(X; \mathcal{D}) - Y\|_2^2 &= \mathbf{E}_{X,Y,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) + \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y\|_2^2 \\
&= \mathbf{E}_{X,Y,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D})\|_2^2 + \mathbf{E}_{X,Y,\mathcal{D}} \|\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,Y,\mathcal{D}} (g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}))^T (\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D})\|_2^2 + \mathbf{E}_{X,Y} \|\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,Y} \mathbf{E}_{\mathcal{D}} ((g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}))^T (\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y) | X, Y) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D})\|_2^2 + \mathbf{E}_{X,Y} \|\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,Y} \mathbf{E}_{\mathcal{D}} (g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) | X, Y)^T (\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D})\|_2^2 + \mathbf{E}_{X,Y} \|\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y\|_2^2 \\
&\quad + 2\mathbf{E}_{X,Y} (\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}))^T (\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y) \\
&= \mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D})\|_2^2 + \mathbf{E}_{X,Y} \|\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y\|_2^2
\end{aligned}$$

Equating the last equation with (10.2) yields

$$\mathbf{E}_{X,Y,\mathcal{D}} \|g(X; \mathcal{D}) - Y\|_2^2 = \underbrace{\mathbf{E}_{X,\mathcal{D}} \|g(X; \mathcal{D}) - \mathbf{E}_{\mathcal{D}}g(X; \mathcal{D})\|_2^2}_{\text{variance}} + \underbrace{\mathbf{E}_{X,Y} \|\mathbf{E}_{\mathcal{D}}g(X; \mathcal{D}) - Y\|_2^2}_{\text{bias} + \text{noise}} \quad (10.3)$$

Therefore in reality, we can only obtain the sum of the bias and noise (not separately) unless we know the quantity of the noise.

Chapter 11

Optimization for Machine Learning

11.1 Gradient method

Suppose that $f : \mathbf{R}^n \rightarrow \mathbf{R}$. An unconstrained optimization problem is

$$\text{minimize } f(x) \tag{11.1}$$

The gradient method is

$$x^{k+1} = x^k - \alpha^k \nabla f(x^k) \tag{11.2}$$

11.2 Stochastic gradient method

Suppose that $f_i : \mathbf{R}^n \rightarrow \mathbf{R}$ for $i = 1, \dots, n$. An unconstrained optimization problem is

$$\text{minimize } \frac{1}{n} \sum_{i=1}^n f_i(x) \tag{11.3}$$

The gradient method is

$$x^{k+1} = x^k - \alpha^k \sum \nabla f(x^k) \tag{11.4}$$

Chapter 12

Bayesian Network

$$\Pr\{X_1, \dots, X_n\} = \prod_{i=1}^n \Pr\{X_i \mid X_1, \dots, X_{i-1}\} = \prod_{i=1}^n \Pr\{X_i \mid \mathbf{parent}(X_i)\} \quad (12.1)$$

Chapter 13

Collaborative Filtering

13.1 Item-based Collaborative Filtering

The purpose of the item-based collaborative filtering is to recommend items using the information obtained from similar items. Since the information from other items help provide better recommendation to customers, it is called an *item-based collaborative filtering*.

We formulate an item-based collaborative filtering as follows. We assume that there are n_C customers and n_I items. We further assume that we have partial information as to the taste of each customer for each item. In typical cases, we have customers can leave reviews with ratings for n items where $n \ll n_I$. When we have hundreds of thousands of items, this makes the data structure extremely sparse. This fact plays a critical role when we calculate similarities among items or customers, or learn the matrix factorization-based collaborative filtering models, e.g., using gradient descent (GD) method or alternating least squares (ALS) method.¹

Now suppose that we have the rating matrix

$$R \in (\mathbf{R} \cup \{\text{NaN}\})^{n_C \times n_I} \quad (13.1)$$

where $R_{ij} \in (\mathbf{R} \cup \{\text{NaN}\})$ denotes *certain score* corresponding to i th customer and j th item and NaN denotes we do not know the values. Note that this certain score can refer to rating of movies, but can mean virtually anything related to customers' taste or activities. For example, it can mean the frequency of customers' clicks on certain menus where the items mean the menu items for this example.

This matrix is sparse, but not in a traditional way. Generally, we say a matrix is sparse if the number of nonzero entries is much less than $n_C \times n_I$. We say R is sparse because there are huge number of entries for which we *do not know* the true values for them. In a general recommendation problem, the purpose is to guess or estimate the true values, e.g., what rating the customer would give if she did, or how frequently a customer would click on the menu item if she knew there exist such menu item.

Now we want to evaluate the similarity (or distance) among items. In many places in the recommendation systems and information retrieval literature, it is shown that applying some transformation or weighting on the values before calculating the similarity measure, e.g., [cosine similarity](#) or [correlation coefficient](#)². There weighting methods and similarity measures will be discussed in later sections.

13.1.1 Rating matrix modeling for menu personalization for mobile shopping app

Here we consider a problem of using a collaborative filtering for efficient menu personalization problem for an mobile shopping app. Suppose that we have some history data of the number of clicks or tabs on each menu item by each customer for certain period. Therefore we can define the matrix in (13.1) where n_C denotes the number of customers, n_I denotes the number of menu items, and R_{ij} denotes the number of clicks or tabs on j th menu item by i th customer with $1 \leq i \leq n_C$ and $1 \leq j \leq n_I$. recommendation system design cases where each component represents the rating

¹These models will be discussed in later sections.

²The correlation coefficient is sometimes called the Pearson correlation coefficient after [Karl Pearson](#), who was an English mathematician and biostatistician.

of an item, these numbers vary depending on the time interval for which we collect the data. So here we model this matrix as a function of time. Moreover, we can have more than one source for the customers' data. Thus we assume multi-source model.

13.1.1.1 Rating matrix modeling

Suppose that $R_s : \mathbf{R} \rightarrow \mathbf{R}^{n_C \times n_I}$ for $s \in \mathcal{S}$ is the rating matrix dependent on time, *i.e.*,

$$R_s(t) \in \mathbf{R}^{n_C \times n_I} \quad (13.2)$$

where \mathcal{S} refers to the set of data sources and $R_s(t)_{ij}$ refers to the number of clicks or tabs for j th item by i th customer which was collected from data source, $s \in \mathcal{S}$. For example, $\mathcal{S} = \{\text{mobile, web}\}$.

For our purpose, we want to have one matrix which represents the customers' activity or behavior with their taste by properly aggregating the history data. Among many possible choices, we choose exponentially decaying weight method together with proper weight on data sources. We define the following aggregate rating matrix

$$\bar{R}(t) = \sum_{s \in \mathcal{S}} \sum_{\tau=0}^{\infty} \gamma^\tau w_s R_s(t - \tau) \quad (13.3)$$

where γ is a positive constant less than 1 and w_s are the weights on each data source such that $\sum_{s \in \mathcal{S}} w_s = 1$.

This aggregate matrix, however, has one problem. Suppose that the 1st customer and the 2nd customer have the very same activity pattern on menu items. However, the 1st customer happens to have actively used the shopping app or website for this week, but the 2nd customer has not used the app or website for the past few days or weeks. Because of the decay factor, the 2nd customer's activity would appear much less than those of the 1st customer even ideally we need to have the very same pattern for those two customers. Therefore we need to normalize the values so that these two customers' preferences are considered the same. For this we normalize each row, so that the sum of each row is always one.

$$\tilde{R}(t) = \mathbf{diag}(\bar{R}(t)\mathbf{1}_{n_I})^{-1} \bar{R}(t) \quad (13.4)$$

where $\mathbf{1}_{n_I} \in \mathbf{R}^{n_I}$ is the vector where every entry is 1 and $\mathbf{diag}(x)$ for $x \in \mathbf{R}^n$ refers to a diagonal matrix whose diagonal entries are the entries of x in the same order. We can readily see that

$$\tilde{R}(t)\mathbf{1}_{n_I} = \mathbf{diag}(\bar{R}(t)\mathbf{1}_{n_I})^{-1} \bar{R}(t)\mathbf{1}_{n_I} = \mathbf{1}_{n_C}, \quad (13.5)$$

i.e., the sum of each row is 1. From this point on, we will remove the subscript for $\mathbf{1}$ unless it can cause confusion as to the dimension of the vector.

One problem of this approach is that $\tilde{R}(t)$ cannot be defined if some row has all zero entries. Even if the sum of every row is nonzero, if the values is very small, *e.g.*, a customer has not shown any activity except for a few clicks or tabs on a handful of menu items, that doesn't mean that the corresponding row represents the customer's preference. Hence, we can augment the values of $\tilde{R}(t)$ by prior distribution of menu items as in the following section.

13.1.2 Value augmentation based on Bayesian MAP

To prevent the divide-by-zero errors from occurring while evaluating (13.4), we consider maximum a posteriori (MAP) estimation for all the rows of $\bar{R}(t)$ in (13.3).

We can think of a problem of filling out each entry in $\bar{R}(t)$ as performing N multinomial experiments on each item $j \in \{1, \dots, n_I\}$ and count the occurrences for each item and fill in $\bar{R}(t)$ for each customer $i \in \{1, \dots, n_C\}$ after normalization.

Now assume that we the prior as Dirichlet-multinomial model, *i.e.*,

$$\text{Dir}(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^{n_I} \theta_k^{\alpha_k-1} I(x=k). \quad (13.6)$$

Then since the likelihood for the multinomial distribution has the form

$$p(\mathcal{D}|\theta) = \prod_{k=1}^{n_I} \theta_k^{N_k}, \quad (13.7)$$

the posterior is

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta) \propto \prod_{k=1}^{n_I} \theta_k^{\alpha_k+N_k-1} = \text{Dir}(\theta|\alpha_1 + N_1, \dots, \alpha_{n_I} + N_{n_I}). \quad (13.8)$$

It can be proved that the maximum a posteriori (MAP) estimate for $\theta_1, \dots, \theta_{n_I}$ is

$$\hat{\theta}_k = \frac{N_k + \alpha_k - 1}{N + \alpha_0 - n_I} \quad (13.9)$$

where $\alpha_0 = \sum_{k=1}^{n_I} \alpha_k$ (K.P. Murphy). By choosing proper values for $\alpha_k > 1$, we can make every entry nonzero in $\bar{R}(t)$.

Adding this information to the rating matrix is equivalent to Bayesian inference since this uses a priori distribution. Therefore, it will play a regularization role in our inference.

13.1.3 Similarity measure among items

Suppose that we have the transformed matrix, $\tilde{R} \in \mathbf{R}^{n_C \times n_I}$. For the case of mobile shopping app menu personalization, $\tilde{R} = \tilde{R}(t)$ for some t . For general cases, \tilde{R} equals to R with all NaNs replaced by 0.

Suppose that $c_1, \dots, c_{n_C} \in \mathbf{R}^{n_I}$ are the row vectors of \tilde{R} and $d_1, \dots, d_{n_I} \in \mathbf{R}^{n_C}$ are the column vectors of \tilde{R} , *i.e.*,

$$\tilde{R} = \begin{bmatrix} c_1^T \\ \vdots \\ c_{n_C}^T \end{bmatrix} = \begin{bmatrix} d_1 & \cdots & d_{n_I} \end{bmatrix} \in \mathbf{R}^{n_C \times n_I}. \quad (13.10)$$

13.1.3.1 Cosine similarity

The cosine similarity measures the cosine of the angle between the two vectors, *i.e.*, it is defined by

$$s(d_i, d_j) = \frac{d_i^T d_j}{\|d_i\|_2 \|d_j\|_2} \quad (13.11)$$

where $x^T y$ denotes the inner product of two vectors x and y , and $\|\cdot\|_2$ denotes the 2-norm of a vector. The [Jensen's inequality](#) guarantees that $-1 \leq s(d_i, d_j) \leq 1$. It can be easily shown that we have $0 \leq s(d_i, d_j) \leq 1$ when every entry in \tilde{R} is nonzero.

13.1.3.2 Cosine similarity when prior distribution is used

When a prior distribution is added to \tilde{R} as in §13.1.2, the sparsity breaks and the matrix becomes a dense matrix. Therefore it seems that we lose huge advantage in computation efforts when calculating the similarities. However, because the added matrix is rank-one matrix, we can still exploit the sparsity and calculate the similarities at almost the same cost as before.

Assume that $\hat{\theta}_k$ in (13.9) has been calculated and is added to \tilde{R} . Thus we have a new rating matrix.

$$\tilde{R}_{\hat{\theta}} = \tilde{R} + \lambda \mathbf{1}_{n_C \times n_I} \mathbf{diag}(\hat{\theta}) = \tilde{R} + \lambda \begin{bmatrix} \hat{\theta}_1 \mathbf{1} & \cdots & \hat{\theta}_{n_I} \mathbf{1} \end{bmatrix} \in \mathbf{R}^{n_C \times n_I} \quad (13.12)$$

where $\hat{\theta} = \begin{bmatrix} \hat{\theta}_1 & \cdots & \hat{\theta}_{n_I} \end{bmatrix}^T \in \mathbf{R}^{n_I}$. Here λ plays a similar role as the coefficient for the regularization when we assume that

$$p(\hat{\theta}) \sim N(0, \lambda I_{n_I}). \quad (13.13)$$

Now let $\tilde{d}_1, \dots, \tilde{d}_{n_I}$ are the column vectors of $\tilde{R}_{\hat{\theta}}$. Then the cosine similarity of i th item and j th item is

$$s(\tilde{d}_i, \tilde{d}_j) = \frac{\tilde{d}_i^T \tilde{d}_j}{\|\tilde{d}_i\|_2 \|\tilde{d}_j\|_2}. \quad (13.14)$$

Now note that

$$\tilde{d}_i^T \tilde{d}_j = (d_i + \lambda \hat{\theta}_i \mathbf{1})^T (d_j + \lambda \hat{\theta}_j \mathbf{1}) = d_i^T d_j + \lambda \hat{\theta}_j \mathbf{1}^T d_i^T + \lambda \hat{\theta}_i \mathbf{1}^T d_j + \lambda^2 \hat{\theta}_i \hat{\theta}_j n_C \quad (13.15)$$

hence

$$\|\tilde{d}_i\|_2^2 = \tilde{d}_i^T \tilde{d}_i = \|d_i\|_2^2 + 2\lambda \hat{\theta}_i \mathbf{1}^T d_i^T + \lambda^2 \hat{\theta}_i^2 n_C \quad (13.16)$$

We can pre-compute $\mathbf{1}^T d_i$ for $i = 1, \dots, n_I$ which takes less than n_R where n_R refers to the number of nonzero entries in \tilde{R} . Hence, the additional computational cost is negligible.

We can also solve this problem at a different angle. In order to calculate the inner projects and the norms in (13.15) and (13.16), we can do the following matrix multiplication.

$$\tilde{R}_{\hat{\theta}}^T \tilde{R}_{\hat{\theta}} = (\tilde{R} + \lambda \mathbf{1}_{n_C \times n_I} \mathbf{diag}(\hat{\theta}))^T (\tilde{R} + \lambda \mathbf{1}_{n_C \times n_I} \mathbf{diag}(\hat{\theta})) \quad (13.17)$$

$$= \tilde{R}^T \tilde{R} + \lambda \tilde{R}^T \mathbf{1}_{n_C \times n_I} \mathbf{diag}(\hat{\theta}) + \lambda \mathbf{diag}(\hat{\theta}) \mathbf{1}_{n_I \times n_C} \tilde{R} \quad (13.18)$$

$$+ \lambda^2 \mathbf{diag}(\hat{\theta}) \mathbf{1}_{n_I \times n_C} \mathbf{1}_{n_C \times n_I} \mathbf{diag}(\hat{\theta}) \quad (13.19)$$

$$= \tilde{R}^T \tilde{R} + \lambda \hat{\theta} \tilde{r}^T + \lambda \tilde{r} \hat{\theta}^T + \lambda^2 n_C \hat{\theta} \hat{\theta}^T \quad (13.20)$$

13.1.3.3 Correlation coefficient similarity

The correlation coefficient is defined by

$$s(d_i, d_j) = \frac{(d_i - \mathbf{1}^T d_i / n_C)^T (d_j - \mathbf{1}^T d_j / n_C)}{\|d_i - \mathbf{1}^T d_i / n_C\|_2 \|d_j - \mathbf{1}^T d_j / n_C\|_2}. \quad (13.21)$$

Again the [Jensen's inequality](#) guarantees that $-1 \leq s(d_i, d_j) \leq 1$, but the fact that every entry in \bar{R} is nonzero does not make this similarity nonnegative.

13.1.4 Data value transformation

When there are popular items across many customers, those values can dominate in the similarity measure evaluation. For example, the aggregate rating matrix, \bar{R} , looks like the following:

$$\begin{bmatrix} \cdots & 100 & \cdots & 130 & \cdots \\ \cdots & 2 & \cdots & 5 & \cdots \\ \cdots & 3 & \cdots & 3 & \cdots \\ \cdots & 10 & \cdots & 1 & \cdots \end{bmatrix} \quad (13.22)$$

$\begin{matrix} \uparrow & & \uparrow \\ d_i & & d_j \end{matrix}$

Note here that the scores in the first row are much larger than the other terms. This can be caused by the fact that some customers are way more active than other customers, *e.g.*, they can listen to some musics many times or uses an shopping app very frequently. Now the cosine similarity and the correlation coefficient similarity between d_i and d_j are

$$0.9956 \text{ and } 0.9951 \quad (13.23)$$

respectively. However, if we remove the first customer and recalculate both similarities, it yields

$$0.4611 \text{ and } -0.9176 \quad (13.24)$$

respectively.

Note that both similarity measures give very different measures. The correlation coefficient similarity, which also tells whether both have positive or negative correlations by its sign, gives opposite signs. This shows how *some* customers activity or behavior can change the item-to-item similarities drastically.

However, this doesn't seem to be right. The item-to-time similarities are supposed to represent the nature of the relation among items, hence should not be decided by a small number of extremely active customers.

For this reason, in many recommendation systems literature, it has been reported that the similarity measures mentioned above do not show good performance. To address this issue, various value transformation methods have been introduced. We will discuss some of these methods in the following sections.

13.1.4.1 TFIDF (or tf-idf)

The [term frequency-inverse document frequency](#) (TFIDF or tf-idf) is a numerical statistic which has been widely used in the field of [information retrieval](#). It was originally designed to reflect how important a word is to a document in a collection or corpus. It is often used as a weighting factor in searches of information retrieval, text mining, and user modeling. The tf-idf value increases proportionally to the number of times a word appears in the document and is offset by the number of documents in the corpus that contain the word, which helps to adjust for the fact that some words appear more frequently in general. Tf-idf is one of the most popular term-weighting schemes today; 83% of text-based recommender systems in digital libraries use tf-idf.

13.1.4.2 Okapi BM25 transformation

Instead, the [Okapi BM25](#) score has shown much better performance when applied to the entries in \tilde{R} before calculating similarity measures. It is defined by

$$\text{bm}_{25}(i, j) = \log \left(\frac{n_C - \|d_j\|_0 + 0.5}{\|d_j\|_0 + 0.5} \right) \frac{\tilde{R}_{ij}(k_1 + 1)}{\tilde{R}_{ij} + k_1 \left(1 - b + b \frac{\|c_i\|_0}{\bar{c}} \right)} \quad (13.25)$$

where $\|\cdot\|_0$ denotes the number of nonzero entries of a vector and k_1 and b are some parameters usually with $k_1 \in [1.2, 2.0]$ and $b \sim 0.75$.

13.1.5 Recommendation based on item similarities

Now finally, we can evaluate new (menu) item scores for each (menu) items for each customer. Here we suggest two methods.

- For customer i and item j , we calculate new recommendation as

$$\hat{R}_{ij} = \frac{\sum_{i=1}^{n_I} s(i, j) \tilde{R}_{i,j}}{\sum_{i=1}^{n_I} s(i, j)} \quad (13.26)$$

where $s(i, j)$ is the cosine similarity between d_i and d_j .

- The second candidate considers the deviation from the average, *i.e.*,

$$\hat{R}_{ij} = \mathbf{1}^T c_i / n_I + \frac{\sum_{i=1}^{n_I} s(i, j) (\tilde{R}_{ij} - \mathbf{1}^T c_i / n_I)}{\sum_{i=1}^{n_I} s(i, j)} \quad (13.27)$$

where $s(i, j)$ is the correlation coefficient similarity between d_i and d_j .

These new values are the ones obtained from the history data for that particular customer together with those for neighboring customers. Hence, this will give better recommendation.

13.2 Collaborative Filtering using Matrix Factorization

Here we discuss the collaborative filtering using matrix factorization, which uses latent factor models. There are two types of latent factors; one for customers and one for items.

The customer latent factors represent customers' taste or tendency. In psychological perspective, even customers themselves may not realize these, but they implicitly express these by their activities. They can be

- the degree of pursuing economical life
- the degree of interest in tableware
- the degree of caring children
- the degree of interest in books
- the degree of putting values on family

On the other hand, the item latent factors represent items' attributes. Again, customers may not realize these, but they tend to click on some menus or items according items attributes and their inclination. They can be

- the probability of leading to inexpensive items
- the probability of leading to tableware
- the probability of leading to items related to children
- the probability of leading to book purchase
- the probability of leading to items related to family

We will discuss how we can utilize these latent factors to design our collaborative filtering below. Note, however, that it is generally impossible to identify the actual latent factors a ML algorithm yields with any of the aforementioned factors. These latent factors are *not predefined or decided a priori*, but are learned or revealed through the values coming out of the algorithms. What's more important is that these methods work very well in practice.

We will also describe different types of problem formulations and different approaches to solve them, *e.g.*, stochastic gradient descent (SGD) method, alternating SGD, and alternating least squares (ALS).

13.2.1 Problem definition and formulations

As before, suppose that there are n_C customers and n_I items. We assume that we have a (sparse) rating matrix, $R \in (\mathbf{R} \cup \{\text{NaN}\})^{n_C \times n_I}$, in (13.1). Sometimes we need to deal with the aggregate rating matrix, $\bar{R}(t) \in (\mathbf{R} \cup \{\text{NaN}\})^{n_C \times n_I}$, in (13.3). To simplify notations, we will refer to both types of rating matrix as R .

Now we assume that there are n_L latent factors. This means that every customer n_L latent factors and every item has n_L latent factors.

Let $x_i \in \mathbf{R}^{n_L}$ be a column vector representing n_L latent factors for the i th customer and let $y_j \in \mathbf{R}^{n_L}$ be a column vector representing n_L latent factors for the j th item with

$$x_i = \begin{bmatrix} x_{i,1} \\ \vdots \\ x_{i,n_L} \end{bmatrix} \in \mathbf{R}^{n_L} \quad (13.28)$$

and

$$y_j = \begin{bmatrix} y_{j,1} \\ \vdots \\ y_{j,n_L} \end{bmatrix} \in \mathbf{R}^{n_L}. \quad (13.29)$$

Then $x_{i,1}, \dots, x_{i,n_L}$ are the n_L tendency or taste factors of the i th customer and let $y_{j,1}, \dots, y_{j,n_L}$ are the n_L attributes of the j th item.

Then we assume that the rating of the j item by the i customer can be inferred (or estimated) by the sum of the corresponding factors, *i.e.*,

$$R_{i,j} \simeq \hat{R}_{i,j} = x_{i,1}y_{j,1} + \dots + x_{i,n_L}y_{j,n_L}. \quad (13.30)$$

This is the most critical assumption in the matrix factorization.

The purpose of the collaborative filtering is to find these latent factors so as to make these inference as accurate as possible, *i.e.*, to solve the following optimization problem

$$\text{minimize} \quad \sum_{1 \leq i \leq n_C, 1 \leq j \leq n_I: R_{i,j} \in \mathbf{R}} l(R_{i,j}, \hat{R}_{i,j}) \quad (13.31)$$

where the $n_L(n_C + n_I)$ optimization variables are $x_{i,k}$ and $y_{j,k}$ for $1 \leq i \leq n_C$, $1 \leq j \leq n_I$, and $1 \leq k \leq n_L$, and $l: \mathbf{R} \times \mathbf{R} \rightarrow \mathbf{R}_+$ is a loss function measuring the distance between the true value and the estimate, hoping that $\hat{R}_{i,j}$ can accurately predict the rating for $1 \leq i \leq n_C$ and $1 \leq j \leq n_I$ where $R_{i,j}$ is not given.

In most cases, we use squared loss for l , so we will also sue the squared loss (except some special cases).

$$l(y_1, y_2) = (y_1 - y_2)^2. \quad (13.32)$$

Now let us come up with more compact notation to describe our problem. Let $X \in \mathbf{R}^{n_C \times n_L}$ be the customer latent factor matrix whose i th row is x_i^T , *i.e.*,

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_{n_C}^T \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,n_L} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,n_L} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n_C,1} & x_{n_C,2} & \cdots & x_{n_C,n_L} \end{bmatrix} \in \mathbf{R}^{n_C \times n_L}. \quad (13.33)$$

Likewise, let $Y \in \mathbf{R}^{n_I \times n_L}$ be the item latent factor matrix whose j th row is y_j^T , *i.e.*,

$$Y = \begin{bmatrix} y_1^T \\ \vdots \\ y_{n_I}^T \end{bmatrix} = \begin{bmatrix} y_{1,1} & y_{1,2} & \cdots & y_{1,n_L} \\ y_{2,1} & y_{2,2} & \cdots & y_{2,n_L} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n_I,1} & y_{n_I,2} & \cdots & y_{n_I,n_L} \end{bmatrix} \in \mathbf{R}^{n_I \times n_L}. \quad (13.34)$$

Note that these two matrices are extremely skinny meaning that we have much more rows than columns in general since n_C could be hundreds of millions and n_I could be hundreds of thousands, but n_L is 100 or so at most.

Now using this notation, we can say

$$\hat{R} = XY^T, \quad (13.35)$$

which is mathematically equivalent to

$$\hat{R}_{i,j} = x_i^T y_j = x_{i,1}y_{j,1} + \cdots + x_{i,n_L}y_{j,n_L}, \quad (13.36)$$

which again is equivalent to (13.30).

Then the optimization problem (13.31) can be rewritten as

$$\text{minimize} \quad \sum_{1 \leq i \leq n_C, 1 \leq j \leq n_I: R_{i,j} \in \mathbf{R}} l(R_{i,j}, x_i^T y_j) \quad (13.37)$$

where the optimization variables are $X \in \mathbf{R}^{n_C \times n_L}$ and $Y \in \mathbf{R}^{n_I \times n_L}$.

If we use the squared loss for l , then we can write this equation more compactly as follows.

$$\text{minimize} \quad \|R - XY^T\|_{F,R}^2 \quad (13.38)$$

where $\|Z\|_{F,R}$ refers to [Frobenius norm](#) calculated for those entries $Z_{i,j}$ with $R_{i,j} \in \mathbf{R}$.

Note that the number of real variables we need to optimize is the same for all the formulations (13.31), (13.37), and (13.38) (because they are equivalent optimization problems).

13.2.2 Solution methods

13.2.2.1 Matrix factorization via singular value decomposition (SVD)

First we discuss obtaining the latent factors, *i.e.*, X and Y , using sparse singular value decomposition (SVD).

For any rank- k matrix $A \in \mathbf{R}^{m \times n}$, we always have a SVD

$$A = U\Sigma V^T \quad (13.39)$$

where $U \in \mathbf{R}^{m \times k}$, $V \in \mathbf{R}^{n \times k}$, and $\Sigma \in \mathbf{R}^{k \times k}$. (Refer to §13.3.1.1.) The basic idea of obtaining the latent factors using SVD is to apply SVD to the rating matrix.

However, the matrix $R \in (\mathbf{R} \cup \{\text{NaN}\})^{n_C \times n_I}$ can have many unknown values, hence we cannot apply SVD directly. Instead, we employ a certain type of missing data imputation to replace the unknowns with real values, then apply SVD to the matrix. There can be many options for the missing data imputation. Here we list some of them.

- replacing every unknown with zeros
- replacing unknowns with global item average ratings
- replacing unknowns with global customer average ratings
- use some model based methods (*e.g.*, the item-based collaborative filtering itself can be considered as such a method)

Now suppose that \tilde{R} is R with missing values replaced by proper real values. We compute the largest n_L singular values with corresponding singular vectors. Then we can obtain the approximation of \tilde{R} by

$$\tilde{R} \simeq U_{n_L} \Sigma_{n_L} V_{n_L}^T \quad (13.40)$$

where

$$U_{n_L} = \begin{bmatrix} u_1 & \cdots & u_{n_L} \end{bmatrix} \in \mathbf{R}^{n_C \times n_L} \quad (13.41)$$

$$V_{n_L} = \begin{bmatrix} v_1 & \cdots & v_{n_L} \end{bmatrix} \in \mathbf{R}^{n_I \times n_L} \quad (13.42)$$

and

$$\Sigma_{n_L} = \mathbf{diag}(\sigma_1, \dots, \sigma_{n_L}) \in \mathbf{R}^{n_L \times n_L}. \quad (13.43)$$

Refer to §13.3.1.1 for more details. This approximation is the best approximation to \tilde{R} in Frobenius norm sense. (Refer to §13.3.1.2 for further details.)

Now since (13.40) can be rewritten as

$$\tilde{R} \simeq U_{n_L} \Sigma_{n_L} V_{n_L}^T = U_{n_L} \Sigma_{n_L}^{1/2} \Sigma_{n_L}^{1/2} V_{n_L}^T = (U_{n_L} \Sigma_{n_L}^{1/2}) (V_{n_L} \Sigma_{n_L}^{1/2})^T \quad (13.44)$$

where

$$\Sigma_{n_L}^{1/2} = \mathbf{diag}(\sigma_1^{1/2}, \dots, \sigma_{n_L}^{1/2}) \in \mathbf{R}^{n_L \times n_L} \quad (13.45)$$

we can obtain the customer and item latent factor matrices as follows.

$$X_{\text{svd}} = U_{n_L} \Sigma_{n_L}^{1/2} = \begin{bmatrix} \sigma_1^{1/2} u_1 & \cdots & \sigma_{n_L}^{1/2} u_{n_L} \end{bmatrix} \in \mathbf{R}^{n_C \times n_L} \quad (13.46)$$

$$Y_{\text{svd}} = V_{n_L} \Sigma_{n_L}^{1/2} = \begin{bmatrix} \sigma_1^{1/2} v_1 & \cdots & \sigma_{n_L}^{1/2} v_{n_L} \end{bmatrix} \in \mathbf{R}^{n_I \times n_L} \quad (13.47)$$

13.2.2.2 Matrix factorization via gradient descent (GD) method

One obvious way to obtain the latent factor matrices is to directly apply [gradient descent method](#) to the following optimization problem directly.

$$\text{minimize } f(X, Y) = \sum_{i,j: R_{i,j} \in \mathbf{R}} l(R_{i,j}, x_i^T y_j) \quad (13.48)$$

which is equivalent to (13.37). Here we denote the objective function of the optimization problem by $f: \mathbf{R}^{n_C \times n_L} \times \mathbf{R}^{n_I \times n_L} \rightarrow \mathbf{R}_+$.

In order to apply gradient descent, we need to evaluate the partial derivative of the objective function with respect to

- $x_{i,k}$ for all $1 \leq k \leq n_L$ and all $i \in \{1 \leq i \leq n_C \mid R_{i,j} \in \mathbf{R} \text{ for some } 1 \leq j \leq n_L\}$
- $y_{j,k}$ for all $1 \leq k \leq n_L$ and all $j \in \{1 \leq j \leq n_I \mid R_{i,j} \in \mathbf{R} \text{ for some } 1 \leq i \leq n_C\}$

For simplicity of the equation derivation, let us assume that

$$\{1 \leq i \leq n_C \mid R_{i,j} \in \mathbf{R} \text{ for some } 1 \leq j \leq n_L\} = \{1, \dots, n_C\} \quad (13.49)$$

$$\{1 \leq j \leq n_I \mid R_{i,j} \in \mathbf{R} \text{ for some } 1 \leq i \leq n_C\} = \{1, \dots, n_I\} \quad (13.50)$$

i.e., every item has at least one rating and every customer has at least one rating.

Now the gradient of $f(X, Y)$ with respect to each x_i is

$$\nabla_{x_i} f(X, Y) = \sum_{j: R_{i,j} \in \mathbf{R}} \frac{\partial}{\partial y_2} l(R_{i,j}, x_i^T y_j) y_j \in \mathbf{R}^{n_L} \text{ for } 1 \leq i \leq n_C \quad (13.51)$$

and that with respect to y_j is

$$\nabla_{y_j} f(X, Y) = \sum_{i: R_{i,j} \in \mathbf{R}} \frac{\partial}{\partial y_2} l(R_{i,j}, x_i^T y_j) x_i \in \mathbf{R}^{n_L} \text{ for } 1 \leq j \leq n_I \quad (13.52)$$

If we use the squared loss for l , the optimization problem becomes

$$\text{minimize } f(X, Y) = \|R - XY^T\|_{F,R}^2 \quad (13.53)$$

which is equivalent to (13.38). Then (13.51) and (13.52) imply that the gradients of $f(X, Y)$ can be calculated by

$$\nabla_{x_i} f(X, Y) = - \sum_{j: R_{i,j} \in \mathbf{R}} (R_{i,j} - x_i^T y_j) y_j \in \mathbf{R}^{n_L} \text{ for } 1 \leq i \leq n_C \quad (13.54)$$

and that with respect to y_j is

$$\nabla_{y_j} f(X, Y) = - \sum_{i: R_{i,j} \in \mathbf{R}} (R_{i,j} - x_i^T y_j) x_i \in \mathbf{R}^{n_L} \text{ for } 1 \leq j \leq n_I \quad (13.55)$$

For notational convenience, we stack these vectors to form derivative matrices as below.

If $R \in \mathbf{R}^{n_C \times n_I}$, *i.e.*, there are no missing values, we can write the gradients more compact form, *i.e.*, as the derivatives of X and Y .

$$D_X f(X, Y) = -2(R - XY^T)Y \in \mathbf{R}^{n_C \times n_L} \quad (13.56)$$

and

$$D_Y f(X, Y) = -2(R^T - YX^T)X \in \mathbf{R}^{n_I \times n_L} \quad (13.57)$$

Finally, we describe the gradient descent method.

- choose learning rate strategy $\eta_k > 0$
- choose initial \tilde{X} and \tilde{Y}
- let $X_0 := \tilde{X}$ and $Y_0 := \tilde{Y}$
- let $k := 0$
- for each iteration, update X and Y

$$X_{k+1} = X_k - \eta_k D_X f(X_k, Y_k) \quad (13.58)$$

$$Y_{k+1} = Y_k - \eta_k D_Y f(X_k, Y_k) \quad (13.59)$$

- stops if certain stopping criterion is satisfied
- update $k := k + 1$ and repeat iterations

13.2.2.3 Matrix factorization via alternating gradient descent (GD) method

Instead of updating X and Y simultaneously, we can update X and Y after the other is updated. The algorithm can be described as follows.

- choose learning rate strategy $\eta_k > 0$
- choose initial \tilde{X} and \tilde{Y}
- let $X_0 := \tilde{X}$ and $Y_0 := \tilde{Y}$
- let $k := 0$
- for each iteration, update X and Y

$$X_{k+1} = X_k - \eta_k D_X f(X_k, Y_k) \quad (13.60)$$

$$Y_{k+1} = Y_k - \eta_k D_Y f(X_{k+1}, Y_k) \quad (13.61)$$

- stops if certain stopping criterion is satisfied
- update $k := k + 1$ and repeat iterations

Note the difference in (13.59) and (13.61). This method updates Y with most recent X values.

13.2.2.4 Matrix factorization via stochastic gradient descent (SGD) method

Theoretically both GD and alternating GD converge (to local minima) with proper learning rate strategies. However, the purpose of the collaborative filtering is not minimize (the sum of) errors (or loss function values), but accurately predict missing ratings, *i.e.*, the ratings that a customer has never given, or sometimes accurately predict the ranking of the items. In this case, what we want to achieve is the solution to the following stochastic optimization problem.

$$\text{minimize } \mathbf{E} \sum_{i,j: R_{i,j} = \text{NaN}} l(\tilde{R}_{i,j}, x_i^T y_j) \quad (13.62)$$

where $\mathbf{E}(\cdot)$ denotes the expected value (or mean/average) of a random variable and $\tilde{R}_{i,j}$ denotes the rating that the i th customer would give to j th item in the future.

Unfortunately there is no direct way to solve this problem because this stochastic optimization problem is not (exactly) solvable. Most importantly, we do not have the future data (some of them would be never available). Therefore solving (13.48) is not what we want.

There is also another problem with solving (13.48); computational cost per iteration. The gradient evaluation takes $3n_L n_R$ multiplications and $3n_L n_R$ additions where n_R refers to the number of known ratings. Thus if the density of R is α , the number of multiplications and additions required to evaluate the gradient is $3\alpha n_L n_C n_I$. Therefore, when n_C or n_I or both are huge, the cost for the gradient calculation can be huge.

To resolve these two problems at the same time, we can use stochastic gradient descent (SGD) method with mini-batch methods. The mini-batch method is to use fixed size of training sets for each gradient descent iteration. This can save the computational cost considerably while approximately solving the stochastic optimization problem (13.62).

13.2.2.5 Matrix factorization via alternating least-squares (ALS)

When we use the squared loss, we solve the problem (13.53). The objective function of this problem is

$$f(X, Y) = \|R - XY^T\|_{F,R}^2 \quad (13.63)$$

This is a quadratic function of a bilinear function of X and Y . Unfortunately, it is not a convex function. If it were a convex function, probably we would not need to use (stochastic) gradient descent method because we have much more efficient methods to solve it (*e.g.*, Newton's method). And if it were a convex function, we would be able to obtain the global minimum (even with gradient descent method).

However, if we fix either X or Y , it becomes the convex function of the other. Indeed, the problem (13.53) becomes a least-squares problem when one of X and Y is fixed, which leads to the alternating least-squares algorithm.

- choose initial \tilde{X} and \tilde{Y}
- let $X_0 := \tilde{X}$ and $Y_0 := \tilde{Y}$
- let $k := 0$
- for each iteration, update X and Y

$$X_{k+1} = \underset{X}{\operatorname{argmin}} \|R - XY_k^T\|_{F,R} \quad (13.64)$$

$$Y_{k+1} = \underset{Y}{\operatorname{argmin}} \|R - X_{k+1}Y^T\|_{F,R} \quad (13.65)$$

- stops if certain stopping criterion is satisfied
- update $k := k + 1$ and repeat iterations

The ALS is named so since solving (13.64) or (13.65) is equivalent to solving a least-squares problem. Note that the functions in (13.64) and (13.65) can be separated by x_i s or y_j s, *e.g.*, the function in (13.64) is

$$f_Y(X) = \|R - XY^T\|_{F,R}^2 = \sum_{i=1}^{n_C} \|\tilde{r}_i^T - x_i^T Y^T\|_{F, \tilde{r}_i^T}^2 = \sum_{i=1}^{n_C} \|\tilde{r}_i - Y x_i\|_{F, \tilde{r}_i}^2 \quad (13.66)$$

where $\tilde{r}_i \in \mathbf{R}^{n_I}$ are the row vectors of R , thus solving (13.64) is equivalent to solving n_C uncorrelated problems separately. This is another great advantage of ALS since we can use parallelism to save the training time considerably. For example, if we have N processing units, the training time per iteration would be equivalent to that for solving each (small) least-squares n_C/N times. Note that we cannot separate the objective function in (13.63) when we consider X and Y simultaneously since each of n_R terms in (13.63) are intertwined through x_i s and y_j s.

When $R \in \mathbf{R}^{n_C \times n_I}$, it can be easily shown that

$$x_i^* = Y^\dagger \tilde{r}_i \quad (13.67)$$

where

$$Y^\dagger = (Y^T Y)^{-1} Y^T \quad (13.68)$$

is the pseudo-inverse of Y . Thus (13.64) becomes

$$X_{k+1} = RY_k^{\dagger T} = RY_k (Y_k^T Y_k)^{-1}. \quad (13.69)$$

Likewise, (13.65) becomes

$$Y_{k+1} = R^T X_{k+1}^{\dagger T} = R^T X_{k+1} (X_{k+1}^T X_{k+1})^{-1}. \quad (13.70)$$

Note that in practice, we do not evaluate the pseudo-inverse as in (13.68) because it causes catastrophic numerical instability especially when the matrix is huge (as in most recommendation system cases). Instead, we use [QR decomposition](#), *i.e.*, if $Y = QR$ is the QR decomposition of Y , $Y^{\dagger} = R^{-1}Q^T$.

We now discuss the interpretation of (13.69) and (13.70). Let X^* and Y^* be the optimal solutions for the problem (13.38). Then these should satisfy (13.64) and (13.65), *i.e.*,

$$X^* = \underset{X}{\operatorname{argmin}} \|R - XY^{*T}\|_{F,R} \quad (13.71)$$

$$Y^* = \underset{Y}{\operatorname{argmin}} \|R - X^*Y^T\|_{F,R} \quad (13.72)$$

Then (13.69) and (13.70) imply that

$$X^*Y^{*T} = RY^*(Y^{*T}Y^*)^{-1}Y^{*T} = RY^{\operatorname{sim}} \quad (13.73)$$

$$= X^*(X^{*T}X^*)^{-1}X^{*T}R = X^{\operatorname{sim}}R \quad (13.74)$$

where X^{sim} and Y^{sim}

$$X^{\operatorname{sim}} = X^*(X^{*T}X^*)^{-1}X^{*T} \in \mathbf{R}^{n_C \times n_C} \quad (13.75)$$

$$Y^{\operatorname{sim}} = Y^*(Y^{*T}Y^*)^{-1}Y^{*T} \in \mathbf{R}^{n_I \times n_I} \quad (13.76)$$

These yield very interesting interpretations for X^{sim} and Y^{sim} . These equations imply that the optimal prediction for the rating of the j th item by the i th customer is

$$x_i^{*T} y_j^* = \sum_{k=1}^{n_I} Y_{k,j}^{\operatorname{sim}} R_{i,k} \quad (13.77)$$

$$= \sum_{k=1}^{n_C} X_{i,k}^{\operatorname{sim}} R_{k,j} \quad (13.78)$$

The equation (13.77) tells us that the optimal prediction is the weighted sum of the i th customer's ratings where the weights are determined by Y^{sim} . This is *exactly the same as the prediction by item-based collaborative filtering* where item-to-item similarity matrix is given by Y^{sim} . On the other hand, the equation (13.78) tells us that the optimal prediction is the weighted sum of the j th item ratings where the weights are determined by X^{sim} . This is *exactly the same as the prediction by user-based collaborative filtering* where user-to-user similarity matrix is given by X^{sim} (which we have not covered in this document).

13.2.2.6 Weighted matrix factorization via alternating least-squares (ALS)

Here we consider the weighted norm for the objective function.

$$f_W(X, Y) = \|W \bullet (R - XY^T)\|_{F,R}^2 \quad (13.79)$$

where $W \in \mathbf{R}_+^{n_C \times n_I}$ and \bullet denotes the component-wise multiplication.

13.2.3 Collaborative filtering for implicit feedback dataset

- binarized variables

$$p_{i,j} = \begin{cases} 1 & R_{i,j} \in \mathbf{R} \\ 0 & R_{i,j} \notin \mathbf{R} \end{cases} \quad (13.80)$$

- confidence variables

$$c_{i,j} = 1 + \alpha R_{i,j} \quad (13.81)$$

or

$$c_{i,j} = 1 + \alpha \log(1 + R_{i,j}/\epsilon) \quad (13.82)$$

- objective function

$$f(X, Y; \lambda) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_I} c_{i,j} (p_{i,j} - x_i^T y_j)^2 + \lambda_X \|X\|_F^2 + \lambda_Y \|Y\|_F^2 \quad (13.83)$$

where $\lambda = (\lambda_X, \lambda_Y) \in \mathbf{R}_{++}^2$ is the coefficient for the regularization.

- gradients

$$\begin{aligned} \nabla_{x_i} f(X, Y; \lambda) &= -2 \sum_{j=1}^{n_I} c_{i,j} (p_{i,j} - x_i^T y_j) y_j + 2\lambda_X x_i \\ &= 2 \left(\left(\sum_{j=1}^{n_I} c_{i,j} y_j y_j^T + \lambda_X I_{n_L} \right) x_i - \sum_{j=1}^{n_I} c_{i,j} p_{i,j} y_j \right) \\ &= 2 \left((Y^T \mathbf{diag}(\tilde{c}_i) Y + \lambda_X I_{n_L}) x_i - Y^T \mathbf{diag}(\tilde{c}_i) \tilde{p}_i \right) \end{aligned} \quad (13.84)$$

Likewise,

$$\nabla_{y_j} f(X, Y; \lambda) = 2 \left((X^T \mathbf{diag}(c_j) X + \lambda_Y I_{n_L}) y_j - X^T \mathbf{diag}(c_j) p_j \right) \quad (13.85)$$

Here $\tilde{c}_i \in \mathbf{R}^{n_I}$ and $c_j \in \mathbf{R}^{n_C}$ are the i th row vector and the j th column vector of $C \in \mathbf{R}^{n_C \times n_I}$ respectively, and $\tilde{p}_i \in \mathbf{R}^{n_I}$ and $p_j \in \mathbf{R}^{n_C}$ are the i th row vector and the j th column vector of $P \in \mathbf{R}^{n_C \times n_I}$ respectively, *i.e.*,

$$C = \begin{bmatrix} \tilde{c}_1 \\ \tilde{c}_2 \\ \vdots \\ \tilde{c}_{n_C} \end{bmatrix} = \begin{bmatrix} c_1 & c_2 & \cdots & c_{n_C} \end{bmatrix} \in \mathbf{R}^{n_C \times n_I} \quad (13.86)$$

$$P = \begin{bmatrix} \tilde{p}_1 \\ \tilde{p}_2 \\ \vdots \\ \tilde{p}_{n_C} \end{bmatrix} = \begin{bmatrix} p_1 & p_2 & \cdots & p_{n_C} \end{bmatrix} \in \mathbf{R}^{n_C \times n_I} \quad (13.87)$$

13.2.3.1 Regularization coefficient conversion

$$f(X, Y; \lambda) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_I} c_{i,j} (p_{i,j} - x_i^T y_j)^2 + \lambda_X \|X\|_F^2 + \lambda_Y \|Y\|_F^2 \quad (13.88)$$

where $\lambda = (\lambda_X, \lambda_Y) \in \mathbf{R}_+^2$.

Suppose that $(X_\lambda^*, Y_\lambda^*)$ is the optimal solution for the following problem:

$$\text{minimize } f(X, Y; \lambda) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_I} c_{i,j} (p_{i,j} - x_i^T y_j)^2 + \lambda_X \|X\|_F^2 + \lambda_Y \|Y\|_F^2 \quad (13.89)$$

for some λ . Since for any X, Y , and $a \neq 0$,

$$f(aX, (1/a)Y; (1/a^2)\lambda_X, a^2\lambda_Y) = f(X, Y; \lambda_X, \lambda_Y), \quad (13.90)$$

for any X and Y ,

$$\begin{aligned} f(aX_\lambda^*, (1/a)Y_\lambda^*; (1/a^2)\lambda_X, a^2\lambda_Y) &= f(X_\lambda^*, Y_\lambda^*; \lambda_X, \lambda_Y) \\ &\leq f((1/a)X, aY; \lambda_X, \lambda_Y) = f(X, Y; (1/a^2)\lambda_X, a^2\lambda_Y). \end{aligned}$$

Therefore $(aX_\lambda^*, (1/a)Y_\lambda^*)$ is the optimal solution for the following problem:

$$\text{minimize } f(X, Y; (1/a^2)\lambda_X, a^2\lambda_Y) \quad (13.91)$$

Therefore if we obtain an optimal solution for a certain $(\tilde{\lambda}_X, \tilde{\lambda}_Y)$, we have readily obtained optimal solutions for all (λ_X, λ_Y) pairs such that $\lambda_X \lambda_Y = \tilde{\lambda}_X \tilde{\lambda}_Y$.

Therefore solving the problem (13.89) is equivalent to solving the following optimization problem:

$$\text{minimize } f(X, Y; \lambda) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_I} c_{i,j} (p_{i,j} - x_i^T y_j)^2 + \sqrt{\lambda_X \lambda_Y} (\|X\|_F^2 + \|Y\|_F^2) \quad (13.92)$$

Now if assume that we have an optimal regularization coefficient $\lambda \in \mathbf{R}_+$ for the following ML problem:

$$\text{minimize } f(X, Y; \lambda) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_I} c_{i,j} (p_{i,j} - x_i^T y_j)^2 + \lambda (\|X\|_F^2 + \|Y\|_F^2) \quad (13.93)$$

If we consider the normalization of each of the three terms in the objective function, we can formulate the problem as follows.

$$\text{minimize } f(X, Y; \lambda) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_I} c_{i,j} (p_{i,j} - x_i^T y_j)^2 / n_C n_I + \lambda (\|X\|_F^2 / n_C n_L + \|Y\|_F^2 / n_I n_L) \quad (13.94)$$

which is equivalent to

$$\text{minimize } f(X, Y; \lambda) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_I} c_{i,j} (p_{i,j} - x_i^T y_j)^2 + \lambda ((n_I / n_L) \|X\|_F^2 + (n_C / n_L) \|Y\|_F^2) \quad (13.95)$$

which again is equivalent to

$$\text{minimize } f(X, Y; \lambda) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_I} c_{i,j} (p_{i,j} - x_i^T y_j)^2 + (\lambda \sqrt{n_C n_I} / n_L) (\|X\|_F^2 + \|Y\|_F^2) \quad (13.96)$$

Now assume that $\lambda^*(n_C, n_I, n_L)$ is the optimal values for λ in (13.93) (which, for example, can be approximately obtained by hyperparameter optimization with the formulation (13.93)). Then, for some other values (n_C', n_I', n_L') , (13.96) implies that the approximate optimal λ can be found by

$$\begin{aligned} \lambda^*(n_C, n_I, n_L) \sqrt{n_C n_I} / n_L &= \lambda^*(n_C', n_I', n_L') \sqrt{n_C' n_I'} / n_L' \\ \Leftrightarrow \lambda^*(n_C', n_I', n_L') &= \lambda^*(n_C, n_I, n_L) n_L' \sqrt{n_C n_I} / n_L \sqrt{n_C' n_I'} \end{aligned}$$

So for example, if we have $n_I = n_I'$ and $n_L = n_L'$, then

$$\lambda^*(n_C', n_I', n_L') = \sqrt{\frac{n_C}{n_C'}} \lambda^*(n_C, n_I, n_L). \quad (13.97)$$

13.3 Appendix

13.3.1 Linear algebra

13.3.1.1 Singular value decomposition (SVD)

For any rank- k matrix $A \in \mathbf{R}^{m \times n}$, there exist three matrix $U \in \mathbf{R}^{m \times k}$, $\Sigma \in \mathbf{R}^{k \times k}$, and $V \in \mathbf{R}^{n \times k}$, such that

$$A = U\Sigma V^T \quad (13.98)$$

where the columns of U are orthonormal, the columns of V are orthonormal, and Σ is a diagonal matrix with nonincreasing positive diagonal entries, *i.e.*,

$$U^T U = V^T V = I_k \in \mathbf{R}^{k \times k} \quad (13.99)$$

and

$$\Sigma = \mathbf{diag}(\sigma_1, \dots, \sigma_k) = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_k \end{bmatrix} \in \mathbf{R}^{k \times k} \quad (13.100)$$

with

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_k > 0 \quad (13.101)$$

where I_k referst to k -by- k identity matrix.

If we let $u_1, \dots, u_k \in \mathbf{R}^n$ be the k column vectors of U , *i.e.*,

$$U = [u_1 \quad \cdots \quad u_k] \in \mathbf{R}^{n \times k}, \quad (13.102)$$

$U^T U = I_k$ holds if and only if

$$u_i^T u_j = \delta_{i,j} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad (13.103)$$

where $\delta_{i,j}$ denotes the [Kronecker delta](#), which means the length of each u_i is 1 and all of them are orthogonal to each other. Likewise, if we let $v_1, \dots, v_k \in \mathbf{R}^m$ be the k column vectors of V , *i.e.*,

$$V = [v_1 \quad \cdots \quad v_k] \in \mathbf{R}^{m \times k}, \quad (13.104)$$

$V^T V = I_k$ holds if and only if

$$v_i^T v_j = \delta_{i,j} \quad (13.105)$$

which means the length of each v_i is 1 and all of them are orthogonal to each other. Note that

- $\sigma_1, \dots, \sigma_k$ are called *singular values* of A .
- u_1, \dots, u_k are called *left singular vectors* of A .
- v_1, \dots, v_k are called *right singular vectors* of A .

Note also that A can be expressed as

$$A = \sum_{i=1}^k \sigma_i u_i v_i^T, \quad (13.106)$$

i.e., A can be express as a linear combination of k rank-1 matrices.

13.3.1.2 Singular value decomposition as rank- k approximation

Given a matrix $A \in \mathbf{R}^{m \times n}$ where $\mathbf{rank}(A) = k$, consider the following optimization problem with $r \leq k$.

$$\begin{aligned} & \text{minimize} && \|A - B\|_F \\ & \text{subject to} && \mathbf{rank}(B) = r \end{aligned} \quad (13.107)$$

where the optimization variable is $B \in \mathbf{R}^{m \times n}$ and $\|\cdot\|_F$ denotes the [Frobenius norm](#).

It turns out that

$$B^* = U_r \Sigma_r V_r^T = \sum_{i=1}^r \sigma_i u_i v_i^T, \quad (13.108)$$

is an optimal solution for the problem (13.107), *i.e.*,

$$\|A - U_r \Sigma_r V_r^T\|_F \leq \|A - B\|_F \quad (13.109)$$

for every B with $\mathbf{rank}(B) = r$ where $U_r \in \mathbf{R}^{m \times r}$, $V_r \in \mathbf{R}^{n \times r}$, and $\Sigma_r \in \mathbf{R}^{r \times r}$ are defined as

$$U_r = \begin{bmatrix} u_1 & \cdots & u_r \end{bmatrix} \in \mathbf{R}^{m \times r} \quad (13.110)$$

$$V_r = \begin{bmatrix} v_1 & \cdots & v_r \end{bmatrix} \in \mathbf{R}^{n \times r} \quad (13.111)$$

$$\Sigma_r = \mathbf{diag}(\sigma_1, \dots, \sigma_r) \in \mathbf{R}^{r \times r} \quad (13.112)$$

This means the close rank- r matrix to A (in Frobenius norm sense) can be obtained with k largest singular values together with corresponding k left and right singular vectors.

Chapter 14

Time Series Anomaly Detection

14.1 Real-Time Anomaly Detection

14.1.1 Computing Anomaly Likelihood

At every time step, we evaluate the moving sample means and the moving sample standard deviations:

$$\mu(t) = \frac{1}{w} \sum_{i=0}^{w-1} s(t-i) \quad (14.1)$$

$$\sigma(t) = \sqrt{\frac{1}{w-1} \sum_{i=0}^{w-1} (s(t-i) - \mu(t))^2} \quad (14.2)$$

Then compute a recent short term average of raw anomaly scores, and apply a threshold to the Gaussian tail probability (Q-function) to decide whether or not to declare an anomaly:

$$L(t) = 1 - Q\left(\frac{\tilde{\mu}(t) - \mu(t)}{\sigma(t)}\right) \quad (14.3)$$

where the short-term moving sample mean is defined by

$$\tilde{\mu}(t) = \frac{1}{w'} \sum_{i=0}^{w'-1} s(t-i) \quad (14.4)$$

and the Q-function is defined by

$$Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^{\infty} \exp\left(-\frac{\tau^2}{2}\right) d\tau. \quad (14.5)$$

We threshold $L(t)$ and report an anomaly if it is very close to 1:

$$L(t) > 1 - \epsilon. \quad (14.6)$$

To take longer history data into consideration while simultaneously emphasizing recent data more, we can consider the exponentially weighted moving sample mean and standard deviation:

$$\mu(t) = (1 - \gamma) \sum_{i=0}^{\infty} \gamma^i s(t-i) \quad (14.7)$$

$$\sigma(t) = \sqrt{(1 - \gamma) \sum_{i=0}^{\infty} \gamma^i (s(t-i) - \mu(t))^2} \quad (14.8)$$

To combine the advantages of finite window method and exponentially weighted method, we can consider the exponentially weighted moving sample mean and standard deviation with finite window

as follows:

$$\mu(t) = \frac{1}{\sum_{i=0}^{w-1} \gamma^i} \sum_{i=0}^{w-1} \gamma^i s(t-i) \quad (14.9)$$

$$\sigma(t) = \sqrt{\frac{1}{\sum_{i=0}^{w-1} \gamma^i} \sum_{i=0}^{w-1} \gamma^i (s(t-i) - \mu(t))^2} \quad (14.10)$$

Chapter 15

Reinforcement Learning

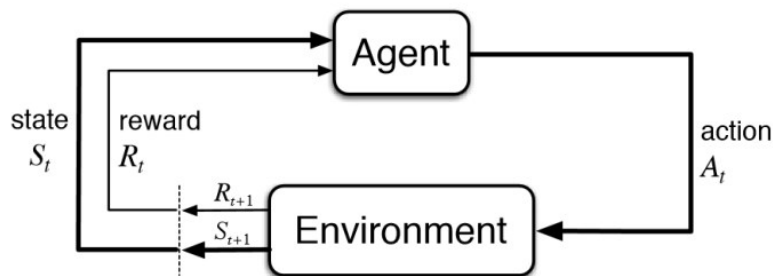


Figure 15.1: The agent-environment interaction in a Markov decision process.

The reinforcement learning is a machine learning where an agent learns how to take actions to achieve a goal by maximizing cumulative reward while interacting with environment. Learning from interaction is a foundational idea underlying nearly all theories of learning and intelligence.

It differs from supervised learning in that labeled input and output pairs need not be presented (and sub-optimal actions need not be explicitly corrected). Instead the focus is finding a balance between exploration of uncharted territory and exploitation of current knowledge. It is much more focused on goal-directed learning from interaction than other approaches to machine learning.

In the following sections, we introduce finite Markov decision process (MDP) and three typical methods to solve reinforcement learning problems. All these methods are called tabular solution methods because they need to store values for all the states or all the state-action pairs (that have been visited).

15.1 Finite Markov decision processes

We introduce the formal problem of finite Markov decision processes (MDPs), which we try to solve. This problem involves evaluative feedback, but also an associative aspect, *i.e.*, choosing different actions in different situations. MDP is a classical formalization of sequential decision making, where actions influence not just immediate rewards, but also subsequent states through those future rewards. Thus MDPs involve delayed reward and the need to trade-off immediate and delayed reward.

Figure 15.1 depicts MDP where an agent interacts with environment. The current state of the environment is known to the agent. With knowledge of state, the agent makes a decision as to which action to take. This action, in turn, will change the state of the environment and the agent will receive a reward at the same time. The agent remembers this reward in some indirect way and uses it for making future decisions.

15.1.1 Markov property

Suppose that the agent is in state S_t takes action A_t at time t . Then the agent receives reward R_{t+1} (from the environment) and the environment transitions to state S_{t+1} . MDP assumes that all these quantities are random variables.

Let \mathcal{S} and \mathcal{A} be the set of all the states and that of all the actions the agent can take respectively.

Now suppose that the environment is in state $S_0 \in \mathcal{S}$ the agent takes action $A_0 \in \mathcal{A}$ at $t = 0$. Then the state of the environment becomes $S_1 \in \mathcal{S}$ giving the agent $R_1 \in \mathbf{R}$ as reward. Suppose that the agent repeat taking actions.

Then we have a sequence of random variables

$$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, S_3, A_3, \dots \quad (15.1)$$

We assume that these random variables satisfy the Markov Property (as assumed by the name) in the following sense.

$$S_{t+1}, R_{t+1} | S_t, A_t, R_t, S_{t-1}, A_{t-1}, R_{t-1}, \dots = S_{t+1}, R_{t+1} | S_t, A_t \quad (15.2)$$

i.e., two random variables, S_{t+1} and R_{t+1} , conditioned on every state, action, and reward before $t + 1$ are the same as those conditioned on S_t and R_t only.

This can be formally expressed using the probability density function (PDF) as follows.

$$p(S_{t+1}, R_{t+1} | S_t, A_t, R_t, S_{t-1}, A_{t-1}, R_{t-1}, \dots) = p(S_{t+1}, R_{t+1} | S_t, A_t). \quad (15.3)$$

This is the reason that the process is called *Markov* decision process.

15.1.2 Policy

The *policy* is defined by the conditional probability of A_t given S_t , *i.e.*,

$$\pi(A|S) = p(A_t | S_t), \quad (15.4)$$

which implies the probability of taking certain action depends only on the current state, not the time. The policy decides which actions the agent takes in each state.

Let Π be the set of all the policies.

15.1.3 Return

The *return* at t is defined by

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \quad (15.5)$$

where $\gamma \in [0, 1]$ is called the *discount factor*. If $\gamma = 0$, the agent is myopic, *i.e.*, it only cares the immediate reward. If $\gamma = 1$, the agent is truly far-sighted, *i.e.*, it cares all the future rewards without discounting. If γ is somewhere between 0 and 1, it considers near-future rewards more importantly than those in far future.

15.1.4 State value function and action value function

The state value function (which is sometimes referred to as just value function) is defined by

$$v_{\pi}(s) = \mathbf{E}_{\pi,p} \{ G_t | S_t = s \} = \mathbf{E}_{\pi,p} \left\{ \sum_{k=0}^{\infty} \gamma^k R_{t+k} \middle| S_t = s \right\}. \quad (15.6)$$

In other words, the state value function is a function of a state representing the expected return the agent will get from the state when following the policy π .

The action value function (which is sometimes referred to as just action function) is defined by

$$q_{\pi}(s, a) = \mathbf{E}_{\pi,p} \{ G_t | S_t = s, A_t = a \} = \mathbf{E}_{\pi,p} \left\{ \sum_{k=0}^{\infty} \gamma^k R_{t+k} \middle| S_t = s, A_t = a \right\}. \quad (15.7)$$

In other words, the action value function is a function of a state and an action representing the expected return the agent will get from the state when the agent takes a certain action and follows the policy π .

As mentioned above, most reinforcement learning algorithms try to maximize either one of these functions, *i.e.*, not maximizing the immediate reward, but the long-term return.

15.2 Bellman equation

[Richard E. Bellman](#), who introduced dynamic programming in 1953, proposed an equation as a necessary condition for optimality associated with dynamic programming, which is called Bellman equation. One of the properties that Markov property implies is that the value functions only depend on the current state (and the action taken) and that the function value is closely related to the function values of the next states. These facts are cleverly used to derive the Bellman equation. Here we introduce two Bellman equations; one for the state value function and the other for action value function.

15.2.1 Bellman equations

To derive Bellman equations, we use some basic statistics facts regarding conditional expectations. (Refer to §15.12.)

Since the definitions of state value function and action value function together with (15.62) imply

$$\begin{aligned} v_{\pi}(s) &= \mathbf{E}_{\pi,p} \{ G_t | S_t = s \} \\ &= \mathbf{E}_{A_t | S_t = s} \mathbf{E}_{\pi,p} \{ G_t | S_t = s, A_t \} \\ &= \sum_a p(A_t = a | S_t = s) \mathbf{E}_{\pi,p} \{ G_t | S_t = s, A_t = a \} \\ &= \sum_a \pi(a | s) \mathbf{E}_{\pi,p} \{ G_t | S_t = s, A_t = a \} \\ &= \sum_a \pi(a | s) q_{\pi}(s, a) \end{aligned}$$

and

$$\begin{aligned}
q_\pi(s, a) &= \mathbf{E}_{\pi, p} \{ G_t | S_t = s, A_t = a \} \\
&= \mathbf{E}_{S_{t+1}, R_{t+1} | S_t = s, A_t = a} \mathbf{E}_{\pi, p} \{ G_t | S_t = s, A_t = a, S_{t+1}, R_{t+1} \} \\
&= \mathbf{E}_{S_{t+1}, R_{t+1} | S_t = s, A_t = a} \mathbf{E}_{\pi, p} \left\{ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s, A_t = a, S_{t+1}, R_{t+1} \right\} \\
&= \mathbf{E}_{S_{t+1}, R_{t+1} | S_t = s, A_t = a} \mathbf{E}_{\pi, p} \left\{ R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} \middle| S_t = s, A_t = a, S_{t+1}, R_{t+1} \right\} \\
&= \sum_{s', r} p_{S_{t+1}, R_{t+1} | S_t, A_t}(s', r | s, a) \mathbf{E}_{\pi, p} \{ R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a, S_{t+1} = s', R_{t+1} = r \} \\
&= \sum_{s', r} p_{S_{t+1}, R_{t+1} | S_t, A_t}(s', r | s, a) (r + \gamma \mathbf{E}_{\pi, p} \{ G_{t+1} | S_t = s, A_t = a, S_{t+1} = s', R_{t+1} = r \}) \\
&= \sum_{s', r} p_{S_{t+1}, R_{t+1} | S_t, A_t}(s', r | s, a) (r + \gamma \mathbf{E}_{\pi, p} \{ G_{t+1} | S_{t+1} = s' \}) \\
&= \sum_{s', r} p_{S_{t+1}, R_{t+1} | S_t, A_t}(s', r | s, a) (r + \gamma v_\pi(s')),
\end{aligned}$$

we have the following two equations relating state value function to action value function and vice versa.

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a). \quad (15.8)$$

$$q_\pi(s, a) = \sum_{s', r} p_{S_{t+1}, R_{t+1} | S_t, A_t}(s', r | s, a) (r + \gamma v_\pi(s')). \quad (15.9)$$

Now (15.8) and (15.9) imply that

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) (r + \gamma v_\pi(s')) \quad (15.10)$$

and

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) (r + \gamma v_\pi(s')) = \sum_{s', r} p(s', r | s, a) \left(r + \gamma \sum_{a'} \pi(a' | s') q_\pi(s', a') \right). \quad (15.11)$$

The equation (15.10) is called *Bellman equation for state value function* and the equation (15.11) is called *Bellman equation for action value function*.

15.2.2 Bellman optimality equations

Now suppose that the policy π_* is the optimal policy. Then we define the *optimal state-value function* as that of π_* , i.e.,

$$v_*(s) = v_{\pi_*}(s) = \max_{\pi \in \Pi} v_\pi(s). \quad (15.12)$$

Likewise, we define the *optimal action-value function* as that of π_* , *i.e.*,

$$q_*(s, a) = q_{\pi_*}(s, a) = \max_{\pi \in \Pi} q_{\pi}(s, a). \quad (15.13)$$

Then (15.8) and (15.9) imply that

$$v_*(s) = v_{\pi_*}(s) = \max_{a \in \mathcal{A}} q_{\pi_*}(s, a) = \max_{a \in \mathcal{A}} \sum_{s', r} p(s', r | s, a) (r + \gamma v_{\pi_*}(s')). \quad (15.14)$$

and

$$q_*(s, a) = q_{\pi_*}(s, a) = \sum_{s', r} p(s', r | s, a) (r + \gamma v_{\pi_*}(s')) = \sum_{s', r} p(s', r | s, a) \left(r + \gamma \max_{a' \in \mathcal{A}} q_{\pi_*}(s', a') \right). \quad (15.15)$$

The equation (15.14) is called *Bellman optimality equation for state value function* and the equation (15.15) is called *Bellman optimality equation for action value function*.

15.3 Dynamic programming

The term dynamic programming (DP) refers to a collection of algorithms that can be used to compute optimal policies given a perfect model of the environment as a Markov decision process (MDP). DP provides an essential foundation for the understanding of the methods presented in the rest of this chapter. All of these methods can be viewed as attempts to achieve much the same effect as DP, only with less computation and without assuming a perfect model of the environment.

The key idea of DP, and of reinforcement learning generally, is the use of value functions to organize and structure the search for good policies.

15.3.1 Policy evaluation (prediction)

We consider how to compute the state-value function v_{π} for an arbitrary policy π . This is called *policy evaluation* in the DP literature. We also refer to it as the *prediction problem*. The existence and uniqueness of v_{π} are guaranteed as long as either $\gamma < 1$ or eventual termination is guaranteed from all states under the policy π .

The policy evaluation algorithm uses the fact that all the state value functions satisfy the Bellman equation for state value function. We use this equation, but in an iterative manner.

$$v_{k+1}(s) \leftarrow \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) (r + \gamma v_k(s')). \quad (15.16)$$

This equation resembles (15.10), but different because now we put subscript k in place of the policy π . Indeed, the sequence v_k can be shown in general to converge to v_{π} as k goes to ∞ the same conditions that guarantee the existence of v_{π} . This algorithm is called *iterative policy evaluation*.

We can consider *in-place* version of this algorithm, *i.e.*, we replace the values for v_k for each state without waiting until we sweep all the states for one iteration. This in-place algorithm also converges to v_{π} . In fact, it usually converges faster. This in-place algorithm is described in Table 15.1.

Inputs: π , MDP
 Algorithm parameters: $\theta > 0$ (small threshold determining accuracy of estimation)

Initialize $V(s) \in \mathbf{R}$ for all $s \in \mathcal{S}$ except that $V(\text{terminal}) = 0$

Loop:
 $\Delta \leftarrow 0$
 For each $s \in \mathcal{S}$:
 $v \leftarrow V(s)$
 $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) (r + \gamma V(s'))$
 $\Delta \leftarrow \max\{\Delta, |v - V(s)|\}$
 until $\Delta < \theta$

Table 15.1: Iterative Policy Evaluation for estimating $V \sim v_\pi$.

15.3.2 Policy iteration

The policy iteration is the iterative process of improving policy as to maximize the value functions. The algorithm is described in Table 15.2.

15.3.3 Value iteration

One drawback to policy iteration is that each of its iterations involves policy evaluation. In fact, the policy evaluation step of policy iteration can be truncated in several ways without losing the convergence guarantees of policy iteration. One important special case is when policy evaluation is stopped after just one sweep (one update of each state). This algorithm is called *value iteration*. It can be written as a particularly simple update operation that combines the policy improvement and truncated policy evaluation steps.

$$v_{k+1}(s) \leftarrow \max_{a \in \mathcal{A}} \sum_{s',r} p(s',r|s,a) (r + \gamma v_k(s')). \quad (15.17)$$

Note that value iteration is obtained simply by turning the Bellman optimality equation for state value function (15.14) into an update rule.

The in-place version of value iteration algorithm is described in Table 15.3.

15.4 Monte Carlo methods

Here we consider learning methods for estimating value functions and discovering optimal policies. Unlike the previous methods, we do not assume complete knowledge of the environment. Monte Carlo (MC) methods require only experience sample sequences of states, actions, and rewards from *actual or simulated interaction with an environment*. Learning from actual experience is striking because it requires no prior knowledge of the environment's dynamics, yet can still attain optimal behavior. *Learning from simulated experience is also powerful*. Although a model is required, the model need only generate sample transitions, *not the complete probability distributions of all possible transitions* that is required for dynamic programming (DP).

<p>Inputs: MDP</p> <p>Algorithm parameters: $\theta > 0$ (small threshold determining accuracy of estimation)</p> <ol style="list-style-type: none"> 1. Initialization $V(s) \in \mathbf{R}$ and $\pi(s) \in \mathcal{A}(s)$ for all $s \in \mathcal{S}$ 2. Policy Evaluation Loop: $\Delta \leftarrow 0$ For each $s \in \mathcal{S}$: $v \leftarrow V(s)$ $V(s) \leftarrow \sum_a \pi(a s) \sum_{s',r} p(s', r s, a) (r + \gamma V(s'))$ $\Delta \leftarrow \max\{\Delta, v - V(s) \}$ until $\Delta < \theta$ 3. Policy Improvement $u \leftarrow \mathbf{true}$ For each $s \in \mathcal{S}$ $b \leftarrow \pi(s)$ $\pi(s) \leftarrow \sum_a \pi(a s) \sum_{s',r} p(s', r s, a) (r + \gamma v_\pi(s'))$ If $b \neq \pi(s)$, then $t \leftarrow \mathbf{false}$ If u, then stop and return $V \sim v_*$ and $\pi \sim \pi_*$; else go to 2
--

Table 15.2: Policy Iteration (using iterative policy evaluation) for estimating $\pi \sim \pi_*$.

<p>Inputs: MDP</p> <p>Algorithm parameters: $\theta > 0$ (small threshold determining accuracy of estimation)</p> <p>Initialize $V(s) \in \mathbf{R}$ for all $s \in \mathcal{S}$ except that $V(\text{terminal}) = 0$</p> <p>Loop: $\Delta \leftarrow 0$ For each $s \in \mathcal{S}$: $v \leftarrow V(s)$ $V(s) \leftarrow \max_{a \in \mathcal{A}(s)} \sum_{s',r} p(s', r s, a) (r + \gamma V(s'))$ $\Delta \leftarrow \max\{\Delta, v - V(s) \}$ until $\Delta < \theta$</p> <p>Output: deterministic policy π such that $\pi(s) = \operatorname{argmax}_{a \in \mathcal{A}(s)} \sum_{s',r} p(s', r s, a) (r + \gamma V(s'))$</p>

Table 15.3: Value Iteration for estimating $\pi \sim \pi_*$.

```

Inputs:  $\pi$ 

Initialize:
   $V(s) \in \mathbf{R}$  for all  $s \in \mathcal{S}$ 
   $R(s) \leftarrow \text{list}()$  for all  $s \in \mathcal{S}$ 

Loop:
  Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
   $G \leftarrow 0$ 
  Loop for each step of episode,  $t = T - 1, T - 2, \dots, 0$ :
     $G \leftarrow \gamma G + R_{t+1}$ 
    If  $S_t \notin \{S_0, S_1, \dots, S_{t-1}\}$ :
       $R(S_t).\text{append}(G)$ 
       $V(S_t) \leftarrow R(S_t).\text{average}()$ 
  Until a certain criterion is satisfied

```

Table 15.4: First-visit MC prediction for estimating $V \sim v_\pi$.

MC methods are ways of solving the reinforcement learning problem based on averaging sample returns. To ensure that well-defined returns are available, here we define Monte Carlo methods only for episodic tasks.

Monte Carlo methods sample and average returns for each state–action pair much like the bandit methods where we sample and average rewards for each action. The main difference is that now there are multiple states, each acting like a different bandit problem (like an associative-search or contextual bandit) and the different bandit problems are interrelated. That is, the return after taking an action in one state depends on the actions taken in later states in the same episode. Because all the action selections are undergoing learning, the problem becomes nonstationary from the point of view of the earlier state.

To handle the nonstationarity, we adapt the idea of general policy iteration (GPI) developed for DP. Whereas there we computed value functions from knowledge of the MDP, here we learn value functions from sample returns with the MDP. The value functions and corresponding policies still interact to attain optimality in essentially the same way (GPI). As in DP, first we consider the prediction problem, then policy improvement, and, finally, the control problem and its solution by GPI. Each of these ideas taken from DP is extended to the Monte Carlo case in which only sample experience is available.

15.4.1 Monte Carlo prediction

An obvious way to estimate it from experience, then, is simply to average the returns observed after visits to that state. There are two Monte Carlo (MC) prediction methods; *first-visit MC method* and *every-visit MC method*. These two MC methods are very similar but have slightly different theoretical properties. First-visit MC has been most widely studied, dating back to the 1940s. Every-visit MC extends more naturally to function approximation and eligibility traces.

Table 15.4 describes the first-visit MC prediction algorithm.

```

Initialize:
   $\pi(s) \in \mathcal{A}(s)$  for all  $s \in \mathcal{S}$ 
   $Q(s, a) \in \mathbf{R}$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ 
   $R(s, a) \leftarrow \text{list}()$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ 

Loop:
  Choose  $S_0 \in \mathcal{S}$ ,  $A_0 \in \mathcal{A}(S_0)$  randomly such that all pairs have probability  $> 0$ 
  Generate an episode from  $S_0, A_0$  following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 
   $G \leftarrow 0$ 
  Loop for each step of episode,  $t + T - 1, T - 2, \dots, 0$ :
     $G \leftarrow \gamma G + R_{t+1}$ 
    If  $S_t \notin \{S_0, S_1, \dots, S_{t-1}\}$ :
       $R(S_t, A_t).\text{append}(G)$ 
       $Q(S_t, A_t) \leftarrow R(S_t, A_t).\text{average}()$ 
       $\pi(S_t) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}(S_t)} Q(S_t, a)$ 
  Until a certain criterion is satisfied

```

Table 15.5: MC ES (exploring starts) for estimating $\pi \sim \pi_*$.

15.4.2 Monte Carlo control

The overall idea is to proceed according to the same pattern as in the dynamic programming, *i.e.*, according to the idea of generalized policy iteration (GPI). In GPI one maintains both an approximate policy and an approximate value function. The value function is repeatedly altered to more closely approximate the value function for the current policy, and the policy is repeatedly improved with respect to the current value function. These two kinds of changes work against each other to some extent, as each creates a moving target for the other, but together they cause both policy and value function to approach optimality.

For Monte Carlo policy evaluation it is natural to alternate between evaluation and improvement on an episode-by-episode basis. After each episode, the observed returns are used for policy evaluation, and then the policy is improved at all the states visited in the episode. A complete simple algorithm along these lines, which is called Monte Carlo ES for Monte Carlo with Exploring Starts, is described in Table 15.5.

15.4.3 Monte Carlo control without exploring starts

How can we avoid the unlikely assumption of exploring starts? The only general way to ensure that all actions are selected infinitely often is for the agent to continue to select them. There are two approaches to ensuring this, resulting in what we call on-policy methods and off-policy methods. On-policy methods attempt to evaluate or improve the policy that is used to make decisions, whereas off-policy methods evaluate or improve a policy different from that used to generate the data. The Monte Carlo ES method developed above is an example of an on-policy method. In this section we show how an on-policy Monte Carlo control method can be designed that does not use the unrealistic assumption of exploring starts.

The on-policy first-visit MC control using ϵ -greedy is described in Table 15.6.

Algorithm parameters: small $\epsilon > 0$

Initialize:

- $\pi(s) \in \mathcal{A}(s)$ for all $s \in \mathcal{S}$
- $Q(s, a) \in \mathbf{R}$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$
- $R(s, a) \leftarrow \text{list}()$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop:

- Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0
- Generate an episode from S_0 , A_0 following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$
- $G \leftarrow 0$
- Loop for each step of episode, $t + T - 1, T - 2, \dots, 0$:
 - $G \leftarrow \gamma G + R_{t+1}$
 - If $S_t \notin \{S_0, S_1, \dots, S_{t-1}\}$:
 - $R(S_t, A_t).\text{append}(G)$
 - $Q(S_t, A_t) \leftarrow R(S_t, A_t).\text{average}()$
 - $A^* \leftarrow \text{argmax}_{a \in \mathcal{A}(S_t)}$
 - For all $a \in \mathcal{A}(S_t)$

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \epsilon + \epsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \epsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

Until a certain criterion is satisfied

Table 15.6: On-policy first-visit MC control (for ϵ -soft policies) for estimating $\pi \sim \pi_*$.

15.4.4 Off-policy prediction via important sampling

XXX

15.4.5 Off-policy Monte Carlo control

XXX

15.5 Temporal-difference learning

Temporal-difference (TD) learning is a combination of Monte Carlo ideas and dynamic programming (DP) ideas. Like Monte Carlo methods, TD methods can learn directly from raw experience without a model of the environment's dynamics. Like DP, TD methods update estimates based in part on other learned estimates, without waiting for a final outcome (they bootstrap). The relationship between TD, DP, and Monte Carlo methods is a recurring theme in the theory of reinforcement learning.

We start by focusing on the policy evaluation or prediction problem, the problem of estimating the value function v_π for a given policy π . For the control problem (finding an optimal policy), DP, TD, and Monte Carlo methods all use some variation of generalized policy iteration (GPI).

15.5.1 TD prediction

Both TD and Monte Carlo methods use experience to solve the prediction problem. A simple every-visit MC method suitable for nonstationary environments is

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t)) = (1 - \alpha)V(S_t) + \alpha G_t. \quad (15.18)$$

TD methods need to wait only until the next time step. At time $t + 1$, they immediately form a target and make a useful update using the observed reward R_{t+1} and the estimate $V(S_{t+1})$. The simplest TD method makes the update

$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)) = (1 - \alpha)V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1})). \quad (15.19)$$

This TD method is called TD(0), or one-step TD, because it is a special case of the TD(λ) and n -step TD methods. Table 15.7 specifies TD(0) completely in procedural form.

The quantity in brackets in the TD(0) update is a sort of error, measuring the difference between the estimated value of S_t and the better estimate $R_{t+1} + \gamma V(S_{t+1})$. This quantity, called the TD error, arises in various forms throughout reinforcement learning. It can be formally defined as follows.

$$\delta_t := R_{t+1} + \gamma V_t(S_{t+1}) - V_t(S_t) \quad (15.20)$$

It is interesting to observe that we can express Monte Carlo error in terms of modified TD errors if we defined the modified TD error as follows.

$$\delta'_t := R_{t+1} + \gamma V_{t+1}(S_{t+1}) - V_t(S_t) \quad (15.21)$$

Inputs: the policy π to be evaluated Algorithm parameters: step size $\alpha \in (0, 1]$ Initialize: $V(s) \in \mathbf{R}$ for all $s \in \mathcal{S}$ except that $V(\text{terminal}) = 0$ Loop for each episode: Initialize S Loop for each step of episode: $A \leftarrow$ action given by π for S Take action A , observe R, S' $V(S) \leftarrow (1 - \alpha)V(S) + \alpha(R + \gamma V(S'))$ $S \leftarrow S'$ until S is terminal Until a certain criterion is satisfied
--

Table 15.7: TD(0) for estimating v_π .

Then the Monte Carlo error is defined by $G_t - V_t(S_t)$ is

$$\begin{aligned}
 G_t - V_t(S_t) &= R_{t+1} + \gamma G_{t+1} - V_t(S_t) \\
 &= R_{t+1} + \gamma (G_{t+1} - V_{t+1}(S_{t+1}) + V_{t+1}(S_{t+1})) - V_t(S_t) \\
 &= R_{t+1} + \gamma V_{t+1}(S_{t+1}) - V_t(S_t) + \gamma (G_{t+1} - V_{t+1}(S_{t+1})) \\
 &= \delta'_t + \gamma (G_{t+1} - V_{t+1}(S_{t+1})) \\
 &= \delta'_t + \gamma \delta'_{t+1} + \gamma^2 (G_{t+2} - V_{t+2}(S_{t+2})) \\
 &= \delta'_t + \gamma \delta'_{t+1} + \gamma^2 \delta'_{t+2} + \cdots + \gamma^{T-t-2} \delta'_{T-2} + \gamma^{T-t-1} (G_{T-1} - V_{T-1}(S_{T-1})) \\
 &= \delta'_t + \gamma \delta'_{t+1} + \gamma^2 \delta'_{t+2} + \cdots + \gamma^{T-t-2} \delta'_{T-2} + \gamma^{T-t-1} (R_T + \gamma V_T(S_T) - V_{T-1}(S_{T-1})) \\
 &= \delta'_t + \gamma \delta'_{t+1} + \gamma^2 \delta'_{t+2} + \cdots + \gamma^{T-t-2} \delta'_{T-2} + \gamma^{T-t-1} \delta'_{T-1} \\
 &= \sum_{k=t}^{T-1} \gamma^{k-t} \delta'_k = \sum_{k=0}^{T-t-1} \gamma^k \delta'_{k+t}
 \end{aligned} \tag{15.22}$$

where the fact that the state-value function for a terminal state, $V_{T-1}(S_T)$, is 0 is used.

This means the Monte Carlo error, *i.e.*, the difference between the return along the path from t to a terminal state of the episode and the state-value function of S_t can be expressed as sum of discounted (modified) one-step TD errors. If we assume that every V_t does not change during the episode, δ_t coincides with δ'_t . Hence (15.22) becomes

$$G_t - V(S_t) = \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k = \sum_{k=0}^{T-t-1} \gamma^k \delta_{k+t}. \tag{15.23}$$

<p>Algorithm parameters: step size $\alpha \in (0, 1]$ and small $\epsilon > 0$</p> <p>Initialize: $Q(s, a) \in \mathbf{R}$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$ except $Q(\text{terminal}, \cdot) = 0$</p> <p>Loop for each episode: Initialize S Choose A from S using policy derived from Q (e.g., ϵ-greedy) Loop for each step of episode: Take action A, observe R, S' Choose A' from S' using policy derived from Q (e.g., ϵ-greedy) $Q(S, A) \leftarrow (1 - \alpha)Q(S, A) + \alpha(R + \gamma Q(S', A'))$ $S \leftarrow S', A \leftarrow A'$, until S is terminal Until a certain criterion is satisfied</p>
--

Table 15.8: Sarsa (on-policy TD control) for estimating $Q \sim q_*$.

15.5.2 Sarsa: on-policy TD Control

As in all on-policy methods, we continually estimate q_π for the behavior policy π , and at the same time change π toward greediness with respect to q_π .

The convergence properties of the Sarsa algorithm depend on the nature of the policy's dependence on Q . For example, one could use ϵ greedy or ϵ -soft policies. Sarsa converges with probability 1 to an optimal policy and action-value function as long as all state-action pairs are visited an infinite number of times and the policy converges in the limit to the greedy policy (which can be arranged, for example, with ϵ -greedy policies by setting $\epsilon = 1/t$).

This algorithm is described in Table 15.8.

15.5.3 Q-learning: off-policy TD control

One of the early breakthroughs in reinforcement learning was the development of an off-policy TD control algorithm known as Q-learning (Watkins, 1989), defined by

$$\begin{aligned}
 Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \alpha \left(R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right) \\
 &= (1 - \alpha)Q(S_t, A_t) + \alpha \left(R_{t+1} + \gamma \max_a Q(S_{t+1}, a) \right).
 \end{aligned} \tag{15.24}$$

The learned action-value function, Q , directly approximates q_* , the optimal action-value function, independent of the policy being followed. This dramatically simplifies the analysis of the algorithm and enabled early convergence proofs. The policy still has an effect in that it determines which state-action pairs are visited and updated. However, all that is required for correct convergence is that all pairs continue to be updated.

Algorithm parameters: step size $\alpha \in (0, 1]$ and small $\epsilon > 0$ Initialize: $Q(s, a) \in \mathbf{R}$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$ except $Q(\text{terminal}, \cdot) = 0$ Loop for each episode: Initialize S Loop for each step of episode: Choose A from S using policy derived from Q (<i>e.g.</i> , ϵ -greedy) Take action A , observe R, S' $Q(S, A) \leftarrow (1 - \alpha)Q(S, A) + \alpha(R + \gamma \max_{a \in \mathcal{A}(S')} Q(S', a))$ $S \leftarrow S'$ until S is terminal Until a certain criterion is satisfied
--

Table 15.9: Q-learning (off-policy TD control) for estimating $\pi \sim \pi_*$.

Under this assumption and a variant of the usual stochastic approximation conditions on the sequence of step-size parameters, Q has been shown to converge with probability 1 to q_* . The Q-learning algorithm is described in Table 15.9.

15.5.4 Maximization bias and double learning

XXX

15.6 *n*-step bootstrapping

There exists another method which unifies the Monte Carlo (MC) methods and the one-step temporal-difference (TD) methods. Neither MC methods nor one-step TD methods are always the best. Here we present *n*-step TD methods that generalize both methods so that one can shift from one to the other smoothly as needed to meet the demands of a particular task. *n*-step methods span a spectrum with MC methods at one end and one-step TD methods at the other. The best methods are often intermediate between the two extremes.

Another way of looking at the benefits of *n*-step methods is that they free one from the tyranny of the time step. With one-step TD methods the same time step determines how often the action can be changed and the time interval over which bootstrapping is done. In many applications one wants to be able to update the action very fast to take into account anything that has changed, but bootstrapping works best if it is over a length of time in which a significant and recognizable state change has occurred. With one-step TD methods, these time intervals are the same, and so a compromise must be made. *n*-step methods enable bootstrapping to occur over multiple steps, freeing us from the tyranny of the single time step.

The idea of *n*-step methods is usually used as an introduction to the algorithmic idea of eligibility traces.

15.6.1 n -step TD prediction

n -step TD prediction is a method lying between Monte Carlo and (one-step) TD method, *i.e.*, TD(0). Consider estimating v_π from sample episodes generated using π . Monte Carlo methods perform an update for each state based on the entire sequence of observed rewards from that state until the end of the episode. The update of one-step TD methods, on the other hand, is based on just the one next reward, bootstrapping from the value of the state one step later as a proxy for the remaining rewards.

One kind of intermediate method, then, would perform an update based on an intermediate number of rewards: more than one, but less than all of them until termination.

The methods that use n -step updates are still TD methods because they still change an earlier estimate based on how it differs from a later estimate. Now the later estimate is not one step later, but n steps later. Methods in which the temporal difference extends over n steps are called n -step TD methods.

Suppose that the process is episodic, *i.e.*, every episode ends or enters a terminal state within finite number of steps. Then the *target* of Monte Carlo update is the return at time step t , *i.e.*,

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{T-t-1} R_T \quad (15.25)$$

where T is the last time step of the episode.

The target of the one-step TD method is the first reward plus the discounted estimated value of the next state, *i.e.*,

$$G_{t:t+1} = R_{t+1} + \gamma V_t(S_{t+1}) \quad (15.26)$$

where $V_t : \mathcal{S} \rightarrow \mathbf{R}$ is the estimate of $v_\pi(S_{t+1})$ at time t . Note that the second term $\gamma V_t(S_{t+1})$ is the estimate for $\gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{T-t-1} R_T$, but using the state-value function estimate of the next state, instead of using future discounted returns. Thus, this is a bootstrapping. Likewise, we can define two-step return as a target for the two-step update.

$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2}) \quad (15.27)$$

where $V_{t+1} : \mathcal{S} \rightarrow \mathbf{R}$ is the estimate of $v_\pi(S_{t+2})$ at time $t+1$. Again here the third term $\gamma^2 V_{t+1}(S_{t+2})$ is the estimate for $\gamma^2 R_{t+3} + \cdots + \gamma^{T-t-1} R_T$ using bootstrapping. In general, we can define the n -step return as the target for the n -step update.

$$G_{t:t+n} = \begin{cases} R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n}) & \text{if } t+n < T \\ G_t & \text{if } t+n \geq T \end{cases} \quad (15.28)$$

for $t \geq 0$. An n -step return can be considered as an approximation to the full return G_t .

We can consider algorithm with the n -step update using this n -step return, *i.e.*,

$$V_{t+n}(S_t) \leftarrow V_{t+n-1}(S_t) + \alpha_{t+n-1}(G_{t:t+n} - V_{t+n-1}(S_t)) = (1 - \alpha_{t+n-1})V_{t+n-1}(S_t) + \alpha_{t+n-1}G_{t:t+n} \quad (15.29)$$

Note that this update cannot be performed before time step $t+n$ because only by then, all the rewards necessary to evaluate (15.28) become available. Also note that (15.29) is reduced to (15.19)

```

Inputs: policy  $\pi$  to be evaluated

Algorithm parameters: step size  $\alpha_t \in (0, 1]$  and  $n \in \mathbf{N}$ 

Initialize:
     $V(s) \in \mathbf{R}$  for all  $s \in \mathcal{S}$  except that  $V(\text{terminal}) = 0$ 

Loop for each episode:
    Initialize and store  $S_0$ 
     $T \leftarrow \infty$ 
    Loop for each  $t = 0, 1, 2, \dots$ :
        If  $t < T$ :
             $A_t \leftarrow$  action given by  $\pi(\cdot|S_t)$ 
            Take action  $A_t$ , observe  $R_{t+1}, S_{t+1}$ 
            If  $S_{t+1}$  is terminal:
                 $T \leftarrow t + 1$ 
             $\tau \leftarrow t - n + 1$ 
            If  $\tau \geq 0$ :
                 $G \leftarrow \sum_{i=\tau+1}^{\min\{\tau+n, T\}} \gamma^{i-\tau-1} R_i$ 
                If  $\tau + n < T$ :
                     $G \leftarrow G + \gamma^n V(S_{\tau+n})$ 
                 $V(S_\tau) \leftarrow (1 - \alpha_t)V(S_\tau) + \alpha_t G$ 
            while  $\tau < T - 1$ 
    Until a certain criterion is satisfied

```

Table 15.10: n -step TD for estimating $V \sim v_\pi$.

```

Inputs: policy  $\pi$  to be evaluated

Algorithm parameters: step size  $\alpha_t \in (0, 1]$  and  $n \in \mathbf{N}$ 

Initialize:
   $V(s) \in \mathbf{R}$  for all  $s \in \mathcal{S}$  except that  $V(\text{terminal}) = 0$ 
   $\mathbf{w\_list} = [1, \gamma, \dots, \gamma^n]$ 

Loop for each episode:
   $\mathbf{S\_list} \leftarrow \text{list}()$ ,  $\mathbf{R\_list} \leftarrow \text{list}()$ 
  Initialize  $S$ 
   $\mathbf{S\_list.append}(S)$ ,  $\text{not\_terminated} \leftarrow \text{True}$ 
  Loop for each  $t = 0, 1, 2, \dots$ :
    If  $\text{not\_terminated}$ :
       $A \leftarrow$  action given by  $\pi(\cdot|S)$ 
      Take action  $A$ , observe  $R, S'$ 
       $\mathbf{R\_list.append}(R)$ 
      If  $S$  is terminal:
         $\text{not\_terminated} \leftarrow \text{False}$ 
      else:
         $\mathbf{S\_list.append}(S')$ 
         $S \leftarrow S'$ 
    If  $t \geq n - 1$ :
       $\mathbf{nR\_list} \leftarrow \mathbf{R\_list}[t - n + 1 : t + 1]$ 
       $G \leftarrow (\mathbf{nR\_list} * \mathbf{w\_list}[: \text{len}(\mathbf{nR\_list})]).\text{sum}()$ 
      If  $t + 1 < \text{len}(\mathbf{S\_list})$ :
         $G \leftarrow G + \mathbf{w\_list}[-1] \times V(\mathbf{S\_list}[t + 1])$ 
       $V(\mathbf{S\_list}[t - n + 1]) \leftarrow (1 - \alpha_t)V(\mathbf{S\_list}[t - n + 1]) + \alpha_t G$ 
    while  $t - n + 1 < \text{len}(\mathbf{R\_list}) - 1$ 
  Until a certain criterion is satisfied

```

Table 15.11: n -step TD for estimating $V \sim v_\pi$ (Pythonic style).

when $n = 1$. Therefore one-step TD method is a special case of n -step TD method. The n -step TD prediction is described in Table 15.10, a Pythonian version of which is described in Table 15.11.

As the Monte Carlo error can be express as the sum of discounted TD errors, the n -step TD error can be expressed as the sum of the discounted one-step TD errors. Note that $G_{t:t+n} = R_{t+1} + \gamma G_{t+1:t+n}$. Thus,

$$\begin{aligned}
G_{t:t+n} - V_t(S_t) &= R_{t+1} + \gamma G_{t+1:t+n} - V_t(S_t) \\
&= R_{t+1} + \gamma (G_{t+1:t+n} + V_{t+1}(S_{t+1}) - V_{t+1}(S_{t+1})) - V_t(S_t) \\
&= R_{t+1} + \gamma V_{t+1}(S_{t+1}) - V_t(S_t) + \gamma (G_{t+1:t+n} - V_{t+1}(S_{t+1})) \\
&= \delta'_t + \gamma (G_{t+1:t+n} - V_{t+1}(S_{t+1})) \\
&= \delta'_t + \gamma \delta'_{t+1} + \gamma^2 (G_{t+2:t+n} - V_{t+2}(S_{t+2})) \\
&= \delta'_t + \gamma \delta'_{t+1} + \gamma^2 \delta'_{t+2} + \cdots + \gamma^{n-1} (G_{t+n-1:t+n} - V_{t+n-1}(S_{t+n-1})) \\
&= \delta'_t + \gamma \delta'_{t+1} + \gamma^2 \delta'_{t+2} + \cdots + \gamma^{n-1} (R_{t+n} + \gamma V_{t+n-1}(S_{t+n}) - V_{t+n-1}(S_{t+n-1})) \\
&= \delta'_t + \gamma \delta'_{t+1} + \gamma^2 \delta'_{t+2} + \cdots + \gamma^{n-1} (R_{t+n} + \gamma V_{t+n}(S_{t+n}) - V_{t+n-1}(S_{t+n-1})) \\
&\quad + \gamma^n (V_{t+n-1}(S_{t+n}) - V_{t+n}(S_{t+n})) \\
&= \delta'_t + \gamma \delta'_{t+1} + \gamma^2 \delta'_{t+2} + \cdots + \gamma^{n-1} \delta'_{t+n-1} + \gamma^n (V_{t+n-1}(S_{t+n}) - V_{t+n}(S_{t+n})) \\
&= \sum_{k=t}^{t+n-1} \gamma^{k-t} \delta'_k + \gamma^n (V_{t+n-1}(S_{t+n}) - V_{t+n}(S_{t+n})) \\
&= \sum_{k=0}^{n-1} \gamma^k \delta'_{k+t} + \gamma^n (V_{t+n-1}(S_{t+n}) - V_{t+n}(S_{t+n}))
\end{aligned}$$

where δ'_t is defined in (15.21). Thus, the n -step TD error is

$$G_{t:t+n} - V_{t+n-1}(S_t) = \sum_{k=0}^{n-1} \gamma^k \delta'_{k+t} + \gamma^n (V_{t+n-1}(S_{t+n}) - V_{t+n}(S_{t+n})) + (V_t(S_t) - V_{t+n-1}(S_t)). \quad (15.30)$$

If V_t does not change during the episode, $\delta_t = \delta'_t$ and $V_{t+n-1}(S_{t+n}) = V_{t+n}(S_{t+n})$, hence

$$G_{t:t+n} - V(S_t) = \sum_{k=t}^{t+n-1} \gamma^{k-t} \delta_k = \sum_{k=0}^{n-1} \gamma^k \delta_{t+k} \quad (15.31)$$

where δ_t is defined in (15.20).

The n -step return uses the value function V_{t+n-1} to correct for the missing rewards beyond R_{t+n} . An important property of n -step returns is that their expectation is guaranteed to be a better estimate of v_π than V_{t+n-1} is, in a worst-state sense. The worst error of the expected n -step return is guaranteed to be less than or equal to γ^n times the worst error under V_{t+n-1} .

$$\max_{s \in S} |\mathbf{E}_\pi \{ G_{t:t+n} | S_t = s \} - v_\pi(s)| \leq \gamma^n \max_s |V_{t+n-1}(s) - v_\pi(s)| \quad (15.32)$$

for all $n \geq 1$.

This is called the error reduction property of n -step returns. Because of the error reduction property, one can show formally that all n -step TD methods converge to the correct predictions under appropriate technical conditions. The n -step TD methods thus form a family of sound methods, with one-step TD methods and Monte Carlo methods as extreme members.

15.6.2 n -step Sarsa

The n -step Sarsa uses the previous n -step temporal difference idea for control. Like one-step Sarsa, we use a behavior policy based on Q functions that the model learns, but update Q functions using the n -step return based on Q function.

$$G_{t:t+n}^Q = \begin{cases} R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n}) & \text{if } t+n < T \\ G_t & \text{if } t+n \geq T \end{cases} \quad (15.33)$$

where $t \geq 0$.

Then, the n -step update using this n -step return is

$$\begin{aligned} Q_{t+n}(S_t, A_t) &= Q_{t+n-1}(S_t, A_t) + \alpha_{t+n-1} (G_{t:t+n}^Q - Q_{t+n-1}(S_t, A_t)) \\ &= (1 - \alpha_{t+n-1}) Q_{t+n-1}(S_t, A_t) + \alpha_{t+n-1} G_{t:t+n}^Q \end{aligned}$$

The algorithm using this update rule is called n -step Sarsa. Table 15.12 describes this algorithm.

As before, the n -step return for Sarsa can be expressed as the sum of TD errors in terms of Q function. Suppose that $t+n < T$. Then

$$\begin{aligned} G_{t:t+n}^Q &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n}) \\ &= \sum_{k=t}^{t+n-1} \gamma^{k-t} R_{k+1} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n}) \\ &= \sum_{k=t}^{t+n-1} \gamma^{k-t} R_{k+1} + \sum_{k=t}^{t+n-1} \gamma^{k-t} (Q_k(S_k, A_k) - Q_k(S_k, A_k)) + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n}) \\ &= \sum_{k=t}^{t+n-1} \gamma^{k-t} R_{k+1} + \sum_{k=t-1}^{t+n-2} \gamma^{k-t+1} Q_{k+1}(S_{k+1}, A_{k+1}) - \sum_{k=t}^{t+n-1} \gamma^{k-t} Q_k(S_k, A_k) \\ &\quad + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n}) \\ &= \sum_{k=t}^{t+n-1} (\gamma^{k-t} R_{k+1} + \gamma^{k-t+1} Q_{k+1}(S_{k+1}, A_{k+1}) - \gamma^{k-t} Q_k(S_k, A_k)) \\ &\quad + Q_t(S_t, A_t) - \gamma^n Q_{t+n}(S_{t+n}, A_{t+n}) + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n}) \\ &= \sum_{k=t}^{t+n-1} \gamma^{k-t} (R_{k+1} + \gamma Q_{k+1}(S_{k+1}, A_{k+1}) - Q_k(S_k, A_k)) \\ &\quad + Q_t(S_t, A_t) + \gamma^n (Q_{t+n-1}(S_{t+n}, A_{t+n}) - Q_{t+n}(S_{t+n}, A_{t+n})) \end{aligned}$$


```

Algorithm parameters: step size  $\alpha_t \in (0, 1]$ , small  $\epsilon > 0$ , and  $n \in \mathbf{N}$ 

Initialize:
     $Q(s, a) \in \mathbf{R}$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$  except that  $Q(\text{terminal}, \cdot) = 0$ 

Loop for each episode:
    Initialize and store  $S_0$ 
    Select an action  $A_0$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
     $T \leftarrow \infty$ 
    Loop for each  $t = 0, 1, 2, \dots$ :
        If  $t < T$ :
            Take action  $A_t$ , observe  $R_{t+1}, S_{t+1}$ 
            If  $S_{t+1}$  is terminal:
                 $T \leftarrow t + 1$ 
            else:
                Select an action  $A_{t+1}$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $\tau \leftarrow t - n + 1$ 
        If  $\tau \geq 0$ :
             $G \leftarrow \sum_{i=\tau+1}^{\min\{\tau+n, T\}} \gamma^{i-\tau-1} R_i$ 
            If  $\tau + n < T$ :
                 $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
             $Q(S_\tau, A_\tau) \leftarrow (1 - \alpha_t) Q(S_\tau, A_\tau) + \alpha_t G$ 
        while  $\tau < T - 1$ 
    Until a certain criterion is satisfied

```

Table 15.12: n -step Sarsa for estimating $Q \sim q_*$ or q_π .

Loop:

Select a state, $S \in \mathcal{S}$, and an action, $A \in \mathcal{A}(S)$, at random

Send S, A to a sample model, and obtain a sample next reward, R , and a sample next state, S'

Apply one-step tabular Q-learning to S, A, R, S' :

$$Q(S, A) \leftarrow (1 - \alpha)Q(S, A) + \alpha (R + \gamma \max_{a \in \mathcal{A}(S)} Q(S', a))$$

Until a certain criterion is satisfied

Table 15.13: Random sample one-step tabular Q-learning.

Thus the n -step error for Sarsa becomes

$$\begin{aligned} G_{t:t+n}^Q - Q_{t+n-1}(S_{t+n}, A_{t+n}) &= \sum_{k=t}^{t+n-1} \gamma^{k-t} (R_{k+1} + \gamma Q_{k+1}(S_{k+1}, A_{k+1}) - Q_k(S_k, A_k)) \\ &\quad + (Q_t(S_t, A_t) - Q_{t+n-1}(S_{t+n}, A_{t+n})) + \gamma^n (Q_{t+n-1}(S_{t+n}, A_{t+n}) - Q_{t+n}(S_{t+n}, A_{t+n})) \end{aligned} \quad (15.34)$$

Again, if Q_t does not change during the episode, (15.34) becomes

$$G_{t:t+n}^Q - Q(S_{t+n}, A_{t+n}) = \sum_{k=t}^{t+n-1} \gamma^{k-t} (R_{k+1} + \gamma Q(S_{k+1}, A_{k+1}) - Q(S_k, A_k)) \quad (15.35)$$

15.6.3 n -step off-policy learning

XXX

15.7 Planning and learning with tabular methods

Table 15.13 describes the random sample one-step tabular Q-learning.

15.7.1 Dyna: integrated planning, acting, and learning

The general Dyna architecture is depicted in Figure 15.2. The tabular Dyna-Q algorithm is described in Table 15.14.

15.8 On-policy Prediction with Approximation

XXX

15.9 On-policy Control with Approximation

XXX

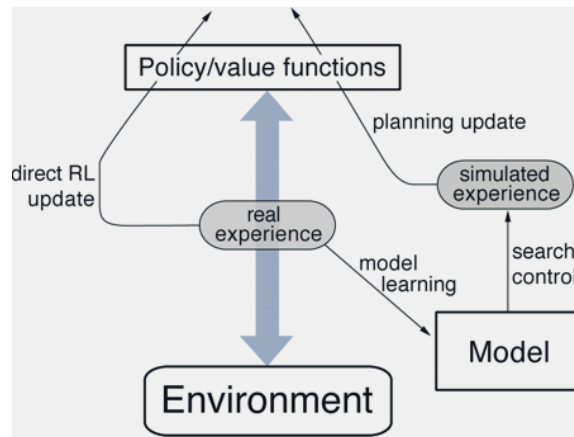


Figure 15.2: The general Dyna Architecture. Real experience, passing back and forth between the environment and the policy, affects policy and value functions in much the same way as does simulated experience generated by the model of the environment.

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(S)$ Loop: $S \leftarrow$ current (nonterminal) state $A \leftarrow \epsilon$ -greedy(S, Q) Take action A ; observe reward R and next state S' Update Q-function: $Q(S, A) \leftarrow (1 - \alpha)Q(S, A) + \alpha (R + \gamma \max_{a \in \mathcal{A}(S)} Q(S', a))$ $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment) Loop repeat Apply one-step tabular Q-learning to S, A, R, S' : $S \leftarrow$ random previously observed state $A \leftarrow$ random action previously taken in S $R, S' \leftarrow Model(S, A)$ Update Q-function: $Q(S, A) \leftarrow (1 - \alpha)Q(S, A) + \alpha (R + \gamma \max_{a \in \mathcal{A}(S)} Q(S', a))$ Until a certain criterion is satisfied
--

Table 15.14: Tabular Dyna-Q.

<p>Algorithm parameters:</p> <p>step size strategy $\alpha_t \in (0, 1]$</p> <p>epsilon strategy: $\epsilon_t > 0$</p> <p>Inputs:</p> <p>Modeling function: $Q : \mathcal{S} \times \mathcal{A} \times \mathbf{R}^n \rightarrow \mathbf{R}$</p> <p>Initialize:</p> <p>Initialize θ arbitrarily (or from the previous learning for transfer learning)</p> <p>Loop for each episode:</p> <p>Initialize S</p> <p>Loop for each step of episode:</p> <p>Choose A from S using policy derived from Q (<i>e.g.</i>, ϵ-greedy)</p> <p>Take action A, observe R, S':</p> $G \leftarrow R + \gamma \max_{a \in \mathcal{A}(S')} Q(S', A; \theta)$ <p>Update using gradient descent:</p> $\theta \leftarrow \theta + \alpha_t (G - Q(S, A; \theta)) \nabla_{\theta} Q(S, A; \theta)$ <p>$S \leftarrow S'$</p> <p>until S is terminal</p> <p>Until a certain criterion is satisfied</p>
--

Table 15.15: Off-policy Q-learning.

15.10 Off-policy Methods with Approximation

15.11 Eligibility Traces

Eligibility traces are one of the basic mechanisms of reinforcement learning. Almost any temporal-difference (TD) method, such as Q-learning or Sarsa, can be combined with eligibility traces to obtain a more general method that may learn more efficiently.

We have already seen one way of unifying TD and Monte Carlo methods: the n -step TD methods §15.6. What eligibility traces offer beyond these is an elegant algorithmic mechanism with significant computational advantages.

The mechanism is a short-term memory vector, the *eligibility trace* $z_t \in \mathbf{R}^d$ that parallels the long-term weight vector $w_t \in \mathbf{R}^d$. The rough idea is that when a component of w_t participates in producing an estimated value, then the corresponding component of z_t is bumped up and then begins to fade away. Learning will then occur in that component of w_t if a nonzero TD error occurs before the trace falls back to zero. The *trace-decay parameter* $\lambda \in [0, 1]$ determines the rate at which the trace falls.

The primary computational advantage of eligibility traces over n -step methods is that only a single trace vector is required rather than a store of the last n feature vectors. Learning also occurs continually and uniformly in time rather than being delayed and then catching up at the end of the episode. In addition learning can occur and affect behavior immediately after a state is encountered rather than being delayed n steps.

15.11.1 The λ -return

In §15.6 we defined an n -step return as the sum of the first n rewards plus the estimated value of the state reached in n steps, each appropriately discounted. The general form of that equation, for any parameterized function approximator, is

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{v}(S_{t+n}, w_{t+n-1}) \quad (15.36)$$

Now we note that a valid update can be done not just toward any n -step return, but toward any average of n -step returns for different ns . An example of such an average can be the average of 3-step return, 5-step return, and 7-step return, *i.e.*,

$$\frac{1}{3}(G_{t:t+3} + G_{t:t+5} + G_{t:t+7}) \quad (15.37)$$

This is one example, but averaging produces a substantial new range of algorithms. For example, we can average across many n -step returns from 1-step return to ∞ -step return to obtain another way of interrelating TD and Monte Carlo methods. On principle, one could even average experience-based updates with DP updates to get a simple combination of experience-based and model-based methods.

An update that averages simpler component updates is called a compound update. The TD(λ) algorithm can be understood as one particular way of averaging n -step updates. This average contains all the n -step updates, each weighted proportionally to λ^n (where $\lambda \in [0, 1]$), and is

normalized (to ensure that the weights sum to 1). The resulting update is toward a return, called the λ -return, defined in its state-based form by

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t:t+n} \quad (15.38)$$

where $(1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} = 1$. If the episode terminates at $t = T$, then $G_{t:t+n} = G_{t:T} = G_t$ for all $n \geq T - t$, thus,

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_{t:t+n} + (1 - \lambda) \sum_{n=T-t}^{\infty} \lambda^{n-1} G_{t:t+n} = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_{t:t+n} + \lambda^{T-t-1} G_t. \quad (15.39)$$

Note that if $\lambda = 1$, G_t^λ becomes the original return G_t . Thus, updating value functions according to λ -return is equivalent to Monte Carlo algorithm. On the other hand, if $\lambda = 0$, G_t^λ becomes $G_{t:t+1}$, hence updating value functions according to λ -return is equivalent to the one-step temporal difference method, *i.e.*, TD(0) method.

Now we define our first learning algorithm based on the λ -return: the *off-line λ -return algorithm*. As an off-line algorithm, it makes no changes to the weight vector during the episode. Then, at the end of the episode, a whole sequence of off-line updates are made according to the semi-gradient method using the λ -return as the target. The update rule of this method is

$$w_{t+1} = w_t + \alpha_t (G_t^\lambda - \hat{v}(S_t, w_t)) \nabla_w \hat{v}(S_t, w_t) \quad (15.40)$$

for $t = 0, \dots, T - 1$. The λ -return gives us an alternative way of moving smoothly between Monte Carlo and one-step TD methods that can be compared with the n -step bootstrapping.

The above approach is what can be called the *theoretical or forward* view of a learning algorithm. For each state visited, we look forward in time to all the future rewards and decide how best to combine them.

15.11.2 TD(λ)

TD(λ) is one of the oldest and most widely used algorithms in reinforcement learning. It was the first algorithm for which a formal relationship was shown between a more theoretical forward view and a more computationally congenial backward view using eligibility traces. Here we will show empirically that it approximates the off-line λ -return algorithm presented in the previous section.

TD(λ) improves over the off-line λ -return algorithm in three ways.

- It updates the weight vector on every step of an episode rather than only at the end, thus its estimates is updated sooner (and may be better sooner).
- Its computations are equally distributed in time (rather than all at the end of the episode).
- It can be applied to continuing problems rather than just to episodic problems.

The semi-gradient version of TD(λ) with function approximation will be shown below.

The *eligibility trace* is a vector $z_t \in \mathbf{R}^d$ where the number of components of z_t is the same as that of w_t . Whereas the *weight vector*, w_t , is a long-term memory accumulating over the lifetime of the system, the eligibility trace is a short-term memory, which typically lasts shorter than the length of an episode. Eligibility traces assist in the learning process; their only consequence is that they affect the weight vector, and then the weight vector determines the estimated value. At the same time, the eligibility traces plays a critical role since that is what propagates future reward to the current one, or equivalently, current reward back to the past states.

In TD(λ), the eligibility trace vector is initialized to zero at the beginning of the episode, is incremented on each time step by the value gradient, and then fades away by $\gamma\lambda$ additively, *i.e.*,

$$\begin{aligned} z_{-1} &\leftarrow 0 \\ z_t &\leftarrow \gamma\lambda z_{t-1} + \nabla_w \hat{v}(S_t, w_t), \quad 0 \leq t \leq T \end{aligned} \quad (15.41)$$

where γ is the *discount rate* and λ is the *trace-decay parameter*.

The eligibility trace keeps track of which components of the weight vector have contributed, positively or negatively, to recent state valuations, where “recent” is defined in terms of $\gamma\lambda$. The trace indicates the eligibility of each component of the weight vector for undergoing learning changes should a reinforcing event occur where the reinforcing events are the moment-by-moment one-step TD errors. The TD error for state-value prediction is

$$\delta_t \leftarrow R_{t+1} + \gamma \hat{v}(S_{t+1}, w_t) - \hat{v}(S_t, w_t) \quad (15.42)$$

Finally, the weight vector is updated on each step proportionally to all three factors; the step size, the TD error, and the eligibility trace, *i.e.*,

$$w_{t+1} \leftarrow w_t + \alpha_t \delta_t z_t. \quad (15.43)$$

Table 15.16 described the *semi-gradient TD(λ) algorithm* for estimating $\hat{v} \sim v_\pi$ using (15.41), (15.42), and (15.43).

15.11.3 Why TD(λ) approximates the off-line λ -return algorithm?

Here we examine why TD(λ) approximates the off-line λ -return algorithm.

In this section, we drop the subscript t from w_t . Because of this, we can only prove that TD(λ) is *approximately* the same as the off-line λ -return algorithm. In fact, they are different.

<p>Inputs:</p> <ul style="list-style-type: none"> π to be evaluated differential function $\hat{v} : \mathcal{S}' \times \mathbf{R}^d \rightarrow \mathbf{R}$ such that $\hat{v}(\text{terminal}, \cdot) = 0$ <p>Algorithm parameters:</p> <ul style="list-style-type: none"> trace-decay parameter, $\lambda \in [0, 1]$, and step size, $\alpha_t > 0$ <p>Initialize:</p> <ul style="list-style-type: none"> Initialize value-function weights w arbitrarily <p>Loop for each episode:</p> <ul style="list-style-type: none"> Initialize S $z \leftarrow 0$ Loop for each step of episode: <ul style="list-style-type: none"> Choose $A \sim \pi(\cdot S)$ Take action A, observe R, S' $z \leftarrow \gamma\lambda z + \nabla_w \hat{v}(S, w)$ $\delta \leftarrow R + \gamma\hat{v}(S', w) - \hat{v}(S, w)$ $w \leftarrow w + \alpha_t \delta z$ $S \leftarrow S'$ until S' is terminal
--

Table 15.16: Semi-gradient TD(λ) algorithm for estimating $\hat{v} \sim v_\pi$.

First, the n -step return at t can be rewritten in terms of that at $t + 1$ as follows.

$$\begin{aligned}
G_{t:t+n} &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{v}(S_{t+n}, w) \\
&= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{n-1} R_{t+n} \\
&\quad + \gamma \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w) \\
&\quad + \gamma(\gamma \hat{v}(S_{t+2}, w) - \hat{v}(S_{t+1}, w)) \\
&\quad + \cdots \\
&\quad + \gamma^{n-1}(\gamma \hat{v}(S_{t+n}, w) - \hat{v}(S_{t+n-1}, w)) \\
&\quad + \hat{v}(S_t, w) \\
&= R_{t+1} + \gamma \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w) \\
&\quad + \gamma(R_{t+2} + \gamma \hat{v}(S_{t+2}, w) - \hat{v}(S_{t+1}, w)) \\
&\quad + \cdots \\
&\quad + \gamma^{n-1}(R_{t+n} + \gamma \hat{v}(S_{t+n}, w) - \hat{v}(S_{t+n-1}, w)) \\
&\quad + \hat{v}(S_t, w),
\end{aligned}$$

thus, (15.42) implies

$$G_{t:t+n} = \delta_t + \gamma \delta_{t+1} + \cdots + \gamma^{n-1} \delta_{t+n-1} + \hat{v}(S_t, w) \quad (15.44)$$

Now (15.44) together with (15.38) implies that

$$\begin{aligned}
G_t^\lambda &= (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t:t+n} \\
&= (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} (\delta_t + \gamma \delta_{t+1} + \cdots + \gamma^{n-1} \delta_{t+n-1} + \hat{v}(S_t, w)) \\
&= (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \left(\sum_{k=0}^{n-1} \gamma^k \delta_{t+k} + \hat{v}(S_t, w) \right) \\
&= (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \left(\sum_{k=0}^{n-1} \gamma^k \delta_{t+k} \right) + (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{v}(S_t, w) \\
&= (1 - \lambda) \sum_{n=1}^{\infty} \sum_{k=0}^{n-1} \lambda^{n-1} \gamma^k \delta_{t+k} + (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{v}(S_t, w) \\
&= (1 - \lambda) \sum_{k=0}^{\infty} \sum_{n=k+1}^{\infty} \lambda^{n-1} \gamma^k \delta_{t+k} + (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{v}(S_t, w) \\
&= (1 - \lambda) \sum_{k=0}^{\infty} \gamma^k \delta_{t+k} \sum_{n=k+1}^{\infty} \lambda^{n-1} + \hat{v}(S_t, w) \\
&= \sum_{k=0}^{\infty} (\gamma \lambda)^k \delta_{t+k} + \hat{v}(S_t, w)
\end{aligned}$$

since $(1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} = 1$, thus

$$G_t^\lambda - \hat{v}(S_t, w) = \sum_{k=0}^{\infty} (\gamma\lambda)^k \delta_{t+k}. \quad (15.45)$$

If an episode ends at $t = T$, then $tderror_t = 0$ for $t \geq T$, hence

$$G_t^\lambda - \hat{v}(S_t, w) = \sum_{k=0}^{T-t-1} (\gamma\lambda)^k \delta_{t+k}. \quad (15.46)$$

Now we derived a formula for z_{t+n} using summation notation from the recursive formula (15.41).

$$\begin{aligned} z_{t+n} &= \gamma\lambda z_{t+n-1} + \nabla_w \hat{v}(S_{t+n}, w) \\ &= \gamma\lambda (\gamma\lambda z_{t+n-2} + \nabla_w \hat{v}(S_{t+n-1}, w)) + \nabla_w \hat{v}(S_{t+n}, w) \\ &= (\gamma\lambda)^2 z_{t+n-2} + (\gamma\lambda) \nabla_w \hat{v}(S_{t+n-1}, w) + \nabla_w \hat{v}(S_{t+n}, w) \\ &\vdots \\ &= (\gamma\lambda)^n \nabla_w \hat{v}(S_t, w) + (\gamma\lambda)^{n-1} \nabla_w \hat{v}(S_{t+1}, w) + \cdots + (\gamma\lambda) \nabla_w \hat{v}(S_{t+n-1}, w) + \nabla_w \hat{v}(S_{t+n}, w), \end{aligned}$$

thus,

$$z_{t+n} = \sum_{k=0}^n (\gamma\lambda)^{n-k} \nabla_w \hat{v}(S_{t+k}, w). \quad (15.47)$$

(We can prove (15.47) mathematically strictly using the mathematical induction.)

We also derive a formula for w_{t+n} using summation notation from the recursive formula (15.43) and (15.42).

$$\begin{aligned} w_{t+n} &= \alpha \delta_{t+n-1} z_{t+n-1} + w_{t+n-1} \\ &= \alpha (\delta_{t+n-1} z_{t+n-1} + \delta_{t+n-2} z_{t+n-2}) + w_{t+n-2} \\ &= \alpha (\delta_{t+n-1} z_{t+n-1} + \delta_{t+n-2} z_{t+n-2} + \cdots + \delta_{t+1} z_{t+1} + \delta_t z_t) + w_t \end{aligned}$$

where we drop the subscript from the step size α_t . Thus,

$$w_{t+n} = \alpha \sum_{j=0}^{n-1} \delta_{t+j} z_{t+j} + w_t. \quad (15.48)$$

Now combining (15.47) and (15.48) yields

$$\begin{aligned}
w_{t+n} - w_t &= \alpha \sum_{j=0}^{n-1} \delta_{t+j} z_{t+j} \\
&= \alpha \sum_{j=0}^{n-1} \delta_{t+j} \sum_{k=0}^j (\gamma \lambda)^{j-k} \nabla_w \hat{v}(S_{t+k}, w). \\
&= \alpha \sum_{j=0}^{n-1} \sum_{k=0}^j (\gamma \lambda)^{j-k} \delta_{t+j} \nabla_w \hat{v}(S_{t+k}, w). \\
&= \alpha \sum_{k=0}^{n-1} \sum_{j=k}^{n-1} (\gamma \lambda)^{j-k} \delta_{t+j} \nabla_w \hat{v}(S_{t+k}, w) \\
&= \alpha \sum_{k=0}^{n-1} \nabla_w \hat{v}(S_{t+k}, w) \sum_{j=k}^{n-1} (\gamma \lambda)^{j-k} \delta_{t+j} \\
&= \alpha \sum_{k=0}^{n-1} \nabla_w \hat{v}(S_{t+k}, w) \sum_{j=0}^{n-k-1} (\gamma \lambda)^j \delta_{t+k+j} \\
&= \alpha \sum_{k=0}^{n-1} \nabla_w \hat{v}(S_{t+k}, w) \sum_{j=0}^{n+t-(t+k)-1} (\gamma \lambda)^j \delta_{(t+k)+j}. \tag{15.49}
\end{aligned}$$

Now if we substitute $T - t$ for n in (15.49), (15.46) implies

$$\begin{aligned}
w_T - w_t &= \alpha \sum_{k=0}^{T-t-1} \nabla_w \hat{v}(S_{t+k}, w) \sum_{j=0}^{T-(t+k)-1} (\gamma \lambda)^j \delta_{(t+k)+j} \\
&= \alpha \sum_{k=0}^{T-t-1} \nabla_w \hat{v}(S_{t+k}, w) (G_{t+k}^\lambda - \hat{v}(S_{t+k}, w)) \\
&= \alpha \sum_{k=0}^{T-t-1} (G_{t+k}^\lambda - \hat{v}(S_{t+k}, w)) \nabla_w \hat{v}(S_{t+k}, w). \\
&= \alpha \sum_{k=t}^{T-1} (G_k^\lambda - \hat{v}(S_k, w)) \nabla_w \hat{v}(S_k, w). \tag{15.50}
\end{aligned}$$

Now observe that the semi-gradient update rule for off-line λ -return algorithm (15.40) implies

that

$$\begin{aligned}
w_T &= \alpha (G_{T-1}^\lambda - \hat{v}(S_{T-1}, w)) \nabla_w \hat{v}(S_{T-1}, w) + w_{T-1} \\
&= \alpha (G_{T-1}^\lambda - \hat{v}(S_{T-1}, w)) \nabla_w \hat{v}(S_{T-1}, w) + \alpha (G_{T-2}^\lambda - \hat{v}(S_{T-2}, w)) \nabla_w \hat{v}(S_{T-2}, w) + w_{T-2} \\
&\vdots \\
&= \alpha (G_{T-1}^\lambda - \hat{v}(S_{T-1}, w)) \nabla_w \hat{v}(S_{T-1}, w) + \cdots + \alpha (G_t^\lambda - \hat{v}(S_t, w)) \nabla_w \hat{v}(S_t, w) + w_t \\
&= \alpha \sum_{k=t}^{T-1} (G_k^\lambda - \hat{v}(S_k, w)) \nabla_w \hat{v}(S_k, w) + w_t
\end{aligned} \tag{15.51}$$

assuming that $\alpha_t = \alpha$ and $w_t = w$ for all t . (We can prove (15.51) in a mathematically strict way using the mathematical induction.) Comparing (15.50) with (15.51) implies that the off-line λ -return algorithm and TD(λ) are equivalent when we drop t from w_t and α_t . Therefore we have shown that TD(λ) approximates the off-line λ -return algorithm.

15.11.4 Sarsa(λ)

15.11.5 Tabular methods using eligibility traces

In the previous sections, we describe TD(λ) for function approximation cases. Since the tabular case is a special case of the function approximation, we can readily derive TD(λ) algorithm for the tabular case. Here we show the process of deriving TD(λ) algorithm for the tabular case and the final algorithm description.

Suppose that the set of all the states, \mathcal{S} , is finite and there is a one-to-one mapping $\phi : \{1, \dots, |\mathcal{S}|\} \rightarrow \mathcal{S}$ which maps each integer index to a state. Let $w \in \mathbf{R}^{|\mathcal{S}|}$ be the weight vector such that $w_i = V(\phi(i))$, i.e., w_i represents the value of the state value function for the i th state, i.e., the state S such with $\phi(i) = S \Leftrightarrow \phi^{-1}(S) = i$. Then we have

$$\hat{v}(S, w) = w_{\phi^{-1}(S)}. \tag{15.52}$$

Now we apply (15.52) to each of the three update rule for TD(λ); (15.41), (15.42), and (15.43). We first derived the gradient of $\hat{v}(S, w)$ with respect to w for each $S \in \mathcal{S}$. Note that

$$\frac{\partial}{\partial w_i} \hat{v}(S, w) = \delta_{\phi^{-1}(S), i} = \begin{cases} 1 & \text{if } \phi(i) = S \\ 0 & \text{otherwise} \end{cases} \tag{15.53}$$

where $\delta_{\cdot, \cdot}$ is the Kronecker delta function defined by

$$\delta_{i,j} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{otherwise.} \end{cases} \tag{15.54}$$

Therefore the gradient is

$$\nabla_w \hat{v}(S, w) = e_{\phi^{-1}(S)} \tag{15.55}$$

where $e_i \in \mathbf{R}^{|\mathcal{S}|}$ is the i th unit vector whose entries are all zero except the i th entry which is 1.

Inputs:
π to be evaluated
Algorithm parameters:
trace-decay parameter, $\lambda \in [0, 1]$, and step size, $\alpha_t > 0$
Initialize:
Initialize $V(S)$ for all $S \in \mathcal{S}$
Loop for each episode:
Initialize S
$z(S) \leftarrow 0$ for all $S \in \mathcal{S}$
Loop for each step of episode:
Choose $A \sim \pi(\cdot S)$
Take action A , observe R, S'
$\delta \leftarrow R + \gamma V(S') - V(S)$
$z(S) \leftarrow z(S) + 1$
For all $\tilde{S} \in \mathcal{S}$:
$V(\tilde{S}) \leftarrow V(\tilde{S}) + \alpha_t \delta z(\tilde{S})$
$z(\tilde{S}) \leftarrow \gamma \lambda z(\tilde{S})$
$S \leftarrow S'$
until S' is terminal

Table 15.17: Tabular TD(λ) algorithm for estimating $\hat{v} \sim v_\pi$.

Now let us apply (15.55) to (15.41), (15.42), and (15.43). First, the eligibility trace update (15.41) becomes

$$\begin{aligned} z_{-1} &\leftarrow \mathbf{0} \in \mathbf{R}^{|\mathcal{S}|} \\ z_t &\leftarrow \gamma \lambda z_{t-1} + e_{\phi^{-1}(S_t)} \quad 0 \leq t \leq T \end{aligned} \tag{15.56}$$

and (15.42) becomes

$$\delta_t \leftarrow R_{t+1} + \gamma w_{\phi^{-1}(S_{t+1})} - w_{\phi^{-1}(S_t)}. \tag{15.57}$$

However, (15.43) remains the same.

$$w_{t+1} \leftarrow w_t + \alpha_t \delta_t z_t. \tag{15.58}$$

Now we note that we would better express this in terms of $V_t : \mathcal{S} \rightarrow \mathbf{R}$ and $z_t : \mathcal{S} \rightarrow \mathbf{R}$. Thus, if we plug (15.56), (15.57), and (15.58) in Table 15.16 using V_t and z_t , we have Table 15.17.

Algorithm parameters:
 trace-decay parameter, $\lambda \in [0, 1]$, and step size, $\alpha_t > 0$

Initialize:
 Initialize $Q(S, A)$ for all $S \in \mathcal{S}$, $A \in \mathcal{A}$

Loop for each episode:
 $z(S, A) \leftarrow 0$ for all $S \in \mathcal{S}$, $A \in \mathcal{A}$
 Initialize S
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 Loop for each step of episode:
 Take action A , observe R , S'
 Choose A' from S' using policy derived from Q (e.g., ϵ -greedy)
 $\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$
 $z(S, A) \leftarrow z(S, A) + 1$
 For all $\tilde{S} \in \mathcal{S}$ and $\tilde{A} \in \mathcal{A}$:
 $Q(\tilde{S}, \tilde{A}) \leftarrow Q(\tilde{S}, \tilde{A}) + \alpha_t \delta z(\tilde{S}, \tilde{A})$
 $z(\tilde{S}, \tilde{A}) \leftarrow \gamma \lambda z(\tilde{S}, \tilde{A})$
 $S \leftarrow S'$, $A \leftarrow A'$
 until S' is terminal

Table 15.18: Tabular SARSA(λ) algorithm for estimating $Q \sim q_*$.

15.12 Appendix: conditional probability and expected value

Suppose that we have a sequence of random variables, X , Y , and Z with supports \mathcal{X} , \mathcal{Y} , and \mathcal{Z} .

Note that the definition of the conditional probability implies that for all $x \in \mathcal{X}$, $y \in \mathcal{Y}$, $z \in \mathcal{Z}$ such that $p_Y(y) \neq 0$ and $p_Z(z) \neq 0$,

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)} \Leftrightarrow p_{X,Y}(x,y) = p_{X|Y}(x|y)p_Y(y). \quad (15.59)$$

$$p(x,y|z) = \frac{p(x,y,z)}{p(z)} = \frac{p(x,y,z)}{p(y,z)} \frac{p(y,z)}{p(z)} = p(x|y,z)p(y|z). \quad (15.60)$$

```

Algorithm parameters:
    trace-decay parameter,  $\lambda \in [0, 1]$ , and step size,  $\alpha_t > 0$ 

Initialize:
    Initialize  $Q(S, A)$  for all  $S \in \mathcal{S}$ ,  $A \in \mathcal{A}$ 

Loop for each episode:
     $z(S, A) \leftarrow 0$  for all  $S \in \mathcal{S}$ ,  $A \in \mathcal{A}$ 
    Initialize  $S$ 
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Loop for each step of episode:
        Take action  $A$ , observe  $R$ ,  $S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $A^* \leftarrow \operatorname{argmin}_{\tilde{A}} Q(S', \tilde{A})$ 
         $\delta \leftarrow R + \gamma Q(S', A^*) - Q(S, A)$ 
         $z(S, A) \leftarrow z(S, A) + 1$ 
        For all  $\tilde{S} \in \mathcal{S}$  and  $\tilde{A} \in \mathcal{A}$ :
             $Q(\tilde{S}, \tilde{A}) \leftarrow Q(\tilde{S}, \tilde{A}) + \alpha_t \delta z(\tilde{S}, \tilde{A})$ 
            If  $\tilde{A} = A^*$ :
                 $z(\tilde{S}, \tilde{A}) \leftarrow \gamma \lambda z(\tilde{S}, \tilde{A})$ 
            else:
                 $z(\tilde{S}, \tilde{A}) \leftarrow 0$ 
         $S \leftarrow S'$ ,  $A \leftarrow A'$ 
    until  $S'$  is terminal

```

Table 15.19: Tabular version of Watkins's $Q(\lambda)$ algorithm for estimating $Q \sim q_*$.

Then (15.59) implies

$$\begin{aligned}
\mathbf{E}(X) &= \int_{\mathcal{X}} xp_X(x)dx \\
&= \int_{\mathcal{X}} x \left(\int_{\mathcal{Y}} p_{X,Y}(x,y)dy \right) dx = \int_{\mathcal{X}} x \left(\int_{\mathcal{Y}} p_{X|Y}(x|y)p_Y(y)dy \right) dx \\
&= \int_{\mathcal{Y}} \int_{\mathcal{X}} xp_{X|Y}(x|y)p_Y(y)dxdy = \int_{\mathcal{Y}} \left(\int_{\mathcal{X}} xp_{X|Y}(x|y)dx \right) p_Y(y)dy \\
&= \int_{\mathcal{Y}} \mathbf{E}(X|Y=y)p_Y(y)dy = \mathbf{E}_Y \mathbf{E}_{X|Y}(X|Y), \\
\mathbf{E}(X) &= \mathbf{E}_{X|Y}(X|Y), \tag{15.61}
\end{aligned}$$

and (15.60) implies

$$\begin{aligned}
\mathbf{E}(X|Z=z) &= \int_{\mathcal{X}} xp_{X|Z}(x|z)dx \\
&= \int_{\mathcal{X}} x \left(\int_{\mathcal{Y}} p_{X,Y|Z}(x,y|z)dy \right) dx = \int_{\mathcal{X}} x \left(\int_{\mathcal{Y}} p_{X|Y,Z}(x|y,z)p_{Y|Z}(y|z)dy \right) dx \\
&= \int_{\mathcal{Y}} \int_{\mathcal{X}} xp_{X|Y,Z}(x|y,z)p_{Y|Z}(y|z)dxdy = \int_{\mathcal{Y}} \left(\int_{\mathcal{X}} xp_{X|Y,Z}(x|y,z)dx \right) p_{Y|Z}(y|z)dy \\
&= \int_{\mathcal{Y}} \mathbf{E}(X|Y=y, Z=z)p_{Y|Z}(y|z)dy = \mathbf{E}_{Y|Z=z} \mathbf{E}(X|Y, Z=z),
\end{aligned}$$

i.e.,

$$\mathbf{E}(X|Z=z) = \mathbf{E}_{Y|Z=z} \mathbf{E}(X|Y, Z=z). \tag{15.62}$$