# "TWO ESSAYS ON DYNAMIC CONTRACTS"

## TESIS PARA OPTAR AL GRADO DE MAGÍSTER EN ECONOMÍA APLICADA

## MEMORIA PARA OPTAR AL TÍTULO DE INGENIERO CIVIL MATEMÁTICO

SOFÍA JOANA MORONI ULLOA

PROFESOR GUÍA:
ALEJANDRO JOFRÉ CÁCERES

MIEMBROS DE LA COMISIÓN:
ANDREA REPETTO LISBOA
NICOLÁS FIGUEROA GONZÁLEZ
FELIPE BALMACEDA MAHNS
HÉCTOR CHADE

SANTIAGO DE CHILE

## TWO ESSAYS ON DYNAMIC CONTRACTS

Esta tesis consta de dos artículos en Teoría de Contratos. El primero es un análisis teórico de las consecuencias de suponer un agente con preferencias aversas a las pérdidas en un modelo de riesgo moral dinámico, usando como base Rogerson (1985). Existe amplia evidencia empírica, tanto en economía como en psicología cognitiva, que demuestra que las preferencias se caracterizan por presentar aversión a las pérdidas. El análisis muestra diferencias relevantes en las predicciones respecto del modelo clásico. En particular, los esquemas óptimos de pagos son no decrecientes pero no estrictamente crecientes en los resultados. En efecto, puede haber segmentos planos en la función de pagos con respecto a los resultados, e incluso puede que en algunos periodos no haya dependencia alguna en éstos. Adicionalmente, se obtiene consecuencias sobre las decisiones de consumo intertemporal del agente. A diferencia de lo que ocurre en Rogerson (1985) se muestra que el agente puede decidir ahorrar, pedir prestado o consumir el pago que recibe si se le otorga acceso al crédito.

Desde el punto de vista metodológico, la aversión a las pérdidas induce una discontinuidad en la utilidad marginal, lo que corresponde a una no diferenciabilidad de la función de utilidad. Luego, para derivar el contrato óptimo se utiliza herramientas de análisis convexo, las que fueron extendidas en este trabajo, introduciendo una nueva regla de la cadena, con el fin de enfrentar el carácter intertemporal del problema. En este trabajo se supone que el punto de referencia se adapta dinámicamente en función del consumo del periodo anterior. Se planifica en trabajo futuro analizar las consecuencias de otras formas de puesta al día de dicha referencia.

El segundo artículo es una aplicación de un modelo agente principal modificado para representar y optimizar acuerdos que se observan entre compañías grandes de servicios IT y sus grandes clientes. Estos contratos son llamados SLA (*Service Level Agreements* o Acuerdos de Nivel de Servicio) y representan un acuerdo entre un proveedor y un consumidor, que explicita objetivos con el fin de garantizar la calidad de servicio. Estos contratos especifican tanto un precio por servicio como penalidades en caso de no cumplimiento. Para la modelación se utiliza conceptos de riesgo moral y selección adversa. El riesgo moral proviene del hecho de que el proveedor, a través de un esfuerzo costoso (inversión, uso de recursos escasos o capital humano, por ejemplo) puede aumentar la calidad del servicio, pero ésta depende además de una componente estocástica. La selección adversa proviene del hecho de que el proveedor puede enfrentar distintos tipos de clientes. La diferencia con el modelo clásico de agente principal redunda en que quien realiza un esfuerzo no verificable es a su vez quien debe determinar el contrato óptimo. En este caso, se incluye una restricción análoga a una de compatibilidad de incentivos, la cual es denominada restricción de credibilidad. Ésta representa el hecho de que el cliente no estaría dispuesto a aceptar un contrato que define un nivel de esfuerzo que no sea óptimo para el proveedor.

Se caracteriza numéricamente el contrato lineal óptimo en un escenario en el que las condiciones están basadas en medidas del tiempo de respuesta de un servicio informático. Adicionalmente, se realiza un análisis de sensibilidad con respecto a parámetros del problema. Se encuentra que los contratos óptimos varían de manera intuitiva al variar parámetros como la aversión al riesgo, la apreciación que tiene el cliente del servicio y otros.

# Abstract

This thesis contains two articles on Contract Theory. The first article is a theoretical analysis of the consequences of assuming that the preferences of the agent are risk averse in the context of a dynamic moral hazard model, using as a benchmark the model in Rogerson (1985). There is considerable empirical evidence, in economics and cognitive psychology that suggests that individual preferences present loss aversion. In this context, the optimal payment scheme has considerable differences with the classical model. In particular, the optimal payment schemes are non-decreasing but not necessarily strictly increasing in outcomes. In fact, the payment scheme can have flat segments that may extend for the entire support of the outcome distribution. In the latter case, the payment scheme is independent of the current period's outcomes. Unlike our benchmark, we find that if given the possibility the agent may decide to save, borrow or consume his allocation under the optimal contract.

Loss aversion induces a discontinuity of marginal utility. Therefore, in order to derive the optimal contract, given the intertemporal properties of the model, it is necessary to use tools from Convex Analysis, and extend them, introducing a new chain rule.

In this work an assumption is that the reference adapts dynamically to the payment received in the last period. Future work could analyze the consequences of other types of update of reference level.

The second article is an application of a modified principal agent model to represent and optimize agreements observed between IT services providers and large clients. These contracts are referred to as SLA (Service Level Agreements) and they are an agreement between a provider and a consumer containing objectives intended to guarantee a quality of service. They specify as well a price for the service and penalties in case of violations. The model of SLA's uses concepts of moral hazard and adverse selection. Moral hazard is identified in the fact that the provider exerts an unobservable costly effort (investment, scarce resources, human capital, for instance) that relates stochastically to the quality of service. The difference between this model and the classical model is that the provider who decides the effort level also determines the contract agreement. In this case, a constraint, analogous an incentive compatibility constraint, is considered, the so-called "credibility constraint". It represents the fact that the client would not be willing to accept an implicit level of effort that is not optimal for the service provider. Adverse selection arises from the different clients the provider may face. The optimal linear contracts are characterized numerically in a scenario in which the terms of the contract are based on measures of response time of an IT service. Additionally, a sensibility analysis is carried out finding that optimal contracts vary as predicted by intuition when modifying parameters such as the risk aversion of the provider or the client, the appreciation the client has of the service and others.

# Acknowledgements

I am deeply indebted to my advisor Prof. Alejandro Jofré whose help, guidance and encouragement helped me during the entire time of research and writing of this thesis.

I also very grateful to Prof. Andrea Repetto for all her helpful comments, suggestions and encouragement that contributed greatly to the completion of this thesis.

I must also thank the members of my commission Felipe Balmaceda, Hector Chade and Nicolas Figueroa for their interest in my work and their useful comments.

Thanks also to the Comisión Nacional de Investigación Científica y Tecnológica (CONICYT) for the support given through the "Beca de Magíster Nacional".

Finally, I'd like to dedicate this work to my parents, my brother and Francisco. I can never thank you enough for your unconditional love and support, always.

# Contents

# List of Figures

# Introduction

This thesis contains two articles, independent of each other, regarding issues of Contract Theory. The first article is a theoretical treatment of the consequences of introducing Loss Averse Preferences with a dynamic update of the reference level in Dynamic Moral Hazard models using as a benchmark the model in Rogerson (1985). Preferences that present loss aversion, as first introduced by Kahneman and Tversky's Prospect Theory (1979), depend on a reference level such that the dislike that consumption below the reference point generates is greater than the elation produced by a gain in the same amount. There is considerable evidence in empirical literature that suggests that references influence individual decisions in economics as well as in cognitive psychology. It is, therefore, relevant to characterize the optimal payment scheme under moral hazard if the Principal is facing an agent whose preferences are reference dependent and this is the objective of the first part of this thesis.

loss aversion induces a discontinuity in marginal utility. Therefore, the optimality conditions that are obtained in classical models are not valid since utility functions are non-differentiable. We derive optimality conditions using convex analysis tools. In summary, the methodology is the following; we show that the program the Principal faces has a concave objective function and the feasible set is convex, therefore, the optimum can be characterized by a subdifferential equals zero condition (Rockafellar, 1974). The computation of the subdifferential in this case is not exempt of mathematical difficulty. It is derived inductively and, in order to so, we must obtain a new chain-rule that applies to the problem at hand. Note that several chain-rule formulas have been described in subdifferential calculus, but none of them applies to the functional forms of our model. Therefore, as a byproduct of this paper we introduce a new chain rule, which allow us to characterize the solutions to an optimization problem of this sort.

The second article is an application of a modified one Principal-Agent Model to a represent and optimize the agreements that a big IT services company contracts with large clients. These contracts are referred to as SLA Contracts. An Service Level Agreement (SLA/Contract) is an agreement between a provider and a consumer which is comprised of Service Level Objectives that guarantee quality of service (such as availability, performance and reliability), a promise of payment and penalties to impose in case the objectives are not met. The study of such contracts has become increasingly important with the increasing use of IT outsourcing procedures, which had reached $56 billion in 2000 and was expected to reach $100 by 2005 (Dermikan et al. (2005)). While the original practice of IT outsourcing contracts involved complicated measures to safeguard the client's interest against the many potential mishaps, a more modern approach has focused on a system of penalties and rewards based on observed quality of service, serving as a monetary compensation that insures the client in case the service is suboptimal (Dermikan et al. (2005)).

In this work we focus on the problem of offering optimal (revenue maximizing) contracts from the Service Providers' (SP) point of view. In particular, we are interested in contracts offered by IT providers, that offer service guarantees in terms of performance, availability, security and reliability constraints. These contracts specify the pricing for the service guarantees and the penalties that are due in case of violations. We model SLA/Contracts using the concepts of Moral Hazard and Adverse Selection.

The Moral Hazard comes from the fact that the provider, through some costly effort (investment, use of scarce resources such as number of CPU's, number of engineer hours, etc.), can increase the quality of the service, but that there is also an additional stochastic component to it. The effort level cannot be monitored by the client, and the actual performance of the system (that the client can observe) is just a noisy signal of effort. In an IT context, better infrastructure on average provides better performance, but some unforeseen incidents (unforeseen demand increase, breakdown of a system, etc.) may still lead to poor quality. Since effort is not observable, the only way to induce a high level of effort is through a compensation system that is "steep" i.e. with higher payments when observed quality is better, or equivalently with penalties if the providers does not meet his end of the deal. Nonetheless, this affects the provider, since she may sometimes be punished for low quality even if the effort put in the process was high. Given her risk aversion, she will demand higher expected payments when the payment system is steeper. The basic trade-off is then set: "steeper" compensation

systems will induce higher effort, but they will shift more risk (in terms of earnings) to the provider, who is risk averse and will charge more for the service. We introduce then the "credibility constraints": a contract must promise a level of effort that is optimal given the penalties imposed in case of non-compliance with the quality level promised. Any other effort level would not be credible and the client would not accept such a contract.

At the same time, the service provider is faced with an adverse selection problem: clients differ in their valuation of the service, in their risk aversion and in other characteristics. Moreover, their particular characteristics are private information, and the service provider only knows the distribution of possible clients. In order to deal with this issue, the service provider must offer different contracts (one for each type of client he can potentially face) and design them in such a way that clients choose the contract that was designed for them. Such constraints (called the self-selection constraints) decrease the revenue an SP may obtain from clients, since they extract an informational rent due to asymmetric information. We construct a general model incorporating both the credibility and self-selection constraints, and allowing for risk averse clients and service provider. Since we are interested in the practical application of such a model to the case of a service provider in the IT sector, we allow for a general shape of the stochastic relation between effort and quality, and proceed to numerically solve for the optimal pricing policy in a particular case. This optimal policy includes different contracts (tailored to be selected by the different types of clients), each one specifying a fixed payment and a bonus based on the quality delivered.

# Chapter 1

# Dynamic Contracts under loss aversion

Alejandro Jofré      Sofia Moroni      Andrea Repetto

## 1.1.   Introduction

The classical principal-agent model is often interpreted as the relationship between an employer, the principal, and an employee, the agent. The agent has private information about the amount of work (or effort) he or she devotes to the task assigned. The outcome obtained, in terms of profits for the principal, relates stochastically to the effort exerted by the agent. In order to induce a desirable level of effort, the employer offers a payment scheme that is contingent on outcomes. Grossman and Hart  (1983) and Hölmstrom  (1979) canonical papers show that under weak assumptions the optimal payment schemes are strictly increasing in outcomes.

In this paper we analyze the consequences of introducing reference dependent preferences that exhibit loss aversion to the canonical model of dynamic moral hazard, as in Rogerson (1985) and Chiappori et al.  (1994). Preferences with loss aversion were presented in Kahneman and Tversky's Prospect Theory (1979).  Under these preferences the dislike that consumption below the reference point generates is greater than the elation produced by a gain in the same amount. There is a large body of literature, in economics and cognitive psychology, that supports that references affect individual decisions (see for example Bateman et al.  (1997)).

De Meza and Webb (2007) first introduced loss aversion in a one period principal-

agent model to find that optimal wage schedules might pay the reference income for different outcomes. In other words, it is possible to observe flat segments at the reference in the optimal payment scheme. The intuition behind this result is that the cost of inducing effort on the reference point through a strictly increasing payment scheme is high due to the sudden decrease in marginal utility for payments below the reference.

In this paper we study a dynamic set up in which the agent's reference is updated endogenously. This type of update is analogous to the one presented in Bowman et al. (1999) and Munro and Sugden (2003). Specifically, we assume that the agent's reference is equal to the previous period's consumption. This new framework modifies some of the predictions of the classical dynamic moral hazard models, as presented in Chiappori et al. (1994), while maintaining others.

We start analyzing the full-commitment case with no access to credit markets, as in Rogerson (1985). In this case, the optimal payment scheme exhibits some characteristics that distinguish it from the classical case. The assumption that each period's reference corresponds to the previous period payment implies that the cost of a payment is not just be the payment itself; it affects the references of later periods and therefore the cost of providing utility to the agent later in the relationship. In our model this effect will tend to lower each period's payment and the slope of the optimal payment scheme, in order to reduce the reference in the following periods, thus, reducing the present value of the cost of inducing incentives.

As in De Meza and Webb (2007), the loss averse preferences of the agent imply that the optimal payment scheme may have a flat segment at the reference. Furthermore, for all periods after the first, the optimal wage schedule must pay the reference for an interval of outcomes. The flat segment may even extend for the whole support of the outcomes distribution. Thus incentives may be optimally provided, not by rewards and punishments that are contingent in the period's results, but by promises of future income. The fact that payment schemes exhibit flat segments in each period implies, under weak assumptions, that there is a positive probability that two consecutive payments are equal. In the canonical model, the same assumptions imply that different outcomes must give different payments to the agent. If the outcomes distribution is continuous no particular payment is given with positive probability. A particular prediction of this model is that in spite of the presence of moral hazard, the optimal schemes of

each period may not depend on the period's outcomes. Only the last period must be contingent on the outcomes realized throughout the relationship. The shape of the optimal contracts obtained in our model may be more easily reconciled with the shape of a variety of contracts that are common in some contexts. For example, contracts that stipulate the same payment for $N - 1$ periods, with a performance evaluation in period $N$ that may derive in a bonus or promotion. Executives that receive options of the firm to be exercised if the firm fares well, in addition to a fixed salary. A "tenure track" type of contract, in which the agent receives a fixed payment for a number of periods in each of which if he excels he is rewarded.

Next we analyze how the agent will allocate resources over time if allowed to save or borrow. We find that the agent can have a somewhat higher incentive to consume the allocation that he was given in a manner consistent with a "status quo bias", as was first described by Samuelson and Zeckhauser (1988). This is because in some situations the agent would experience a loss in utility by both saving and borrowing. loss aversion implies that the marginal utility of savings may not be equal to minus the marginal utility of borrowing and they may be both be negative simultaneously. Also, as long as the agent is consuming his reference in one period or there is positive probability of being paid the reference in the next period, we find a gap between the willingness to lend (or save) and willingness to borrow. The smallest interest rate at which the agent would be willing to accept lending part of his income is strictly higher than the rate he is willing to pay to increase his income. In other words, there is a gap between the price at which he is willing to lend part of his income and the price at which he is willing to borrow. This gap between willingness to accept and willingness to pay has been described in other economical contexts (Bateman et al. (1997)).

The classical result by Rogerson (1985) asserts that in the canonical model, if allowed to save or borrow the agent would want to save after the each period will not necessarily be true in our model. The agent may have incentives to borrow or save after each period depending on the parameters of the problem. This is partly due to the fact that savings and borrowing not only modifies the inter-temporal allocation of consumption but also changes the future references of the agent. In the classical model the optimal payment scheme required the agent to consume more than he would like in order to facilitate the provision of incentives in future periods. In our model, this effect is also present, however, there are other effects at stake that can make the agent want

to borrow or consume his allocation after some periods. First, the fact that the agent has a relatively higher loss in utility towards losses may make it too costly to provide incentives with payments below the reference and thus the optimal payment scheme may pay below the reference for few or no outcomes. Knowing that there is a small likelihood to receive payments below the reference the agent may have incentives to consume his allocation or to borrow. Second, the principal also has incentives to reduce the payment of the agent in each period in order to increase the utility in the following period, again providing incentives to borrow. Finally, for two consecutive periods, if payments are at the reference, the agent faces a high marginal loss in the present that is not compensated by the marginal gain due to savings in the future.

Nevertheless, the agent may also have incentives to save under some circumstances. It is less costly to provide incentives with payments below the reference because of the relatively higher marginal utility; therefore, an optimal payment scheme may give payments below the reference for a large set of outcomes. This would favor a desire of saving after some periods. Also, the agent realizes that a high consumption will decrease his future utility thus giving incentives to save. Finally, as in the classical case the Principal has incentives to provide higher payments at the beginning of the relationship in order to facilitate the provision of incentives in future periods.

The fact that the optimal payment scheme does not always require constrained savings is an important improvement over the canonical model. In the classical case full-commitment requires constrained savings, and lifting the assumption of constrained credit implies that renegotiation proofness is broken. Also, a renegotiation-proof long-term contract with free access to credit cannot provide incentives to exert an effort over the minimum after the first period. Since it is unlikely that a court of law would hinder renegotiation towards a Pareto improving agreement and constraining savings may be implausible in most contexts, the classical theory cannot explain the existence of long-term commitment contracts (see Chiappori et al. (1994)). Therefore, the fact that loss aversion and a dynamic update of the reference does not require to constrain savings in some cases, and is ex-post efficient, might give a rationale for the ubiquity of commitment contracts. For instance, under borrowing constraints a contract that doesn't constrain savings is ex-post efficient. Empirical literature suggests that constrained borrowing is, in fact, present in many different contexts (see for instance Carroll and Kimball (2001)).

Next we analyze under which conditions the optimal contract is over the reference in every period. When this happens the payment scheme presents a ratchet effect. Once the principal gives the agent a payment the following period's scheme must be greater or equal than that payment. From the optimality conditions it is clear that the shape of the optimal payment scheme relates to the marginal costs of the constraints faced by the principal, identified by the corresponding multipliers. We find that for any given marginal cost of the incentive compatibility constraints there is a threshold for the marginal cost of the participation constraint over which payments are always above the reference. Thus, if the participation constraint is very costly for the principal, the optimal payment schemes do not venture into the loss area and incentives will be created either with gains for good outcomes or promises of future income. Later we show that the agent does not have incentives to save under that kind of payment schemes.

We analyze other scenarios under the loss averse preferences of our model. We characterize the optimal sequence of spot contracts and show that will not be memoryless as in the classical case and it will not coincide with the full-commitment optimum in general. It is straightforward to show that the full-commitment optimum is ex-post efficient and therefore renegotiation proof. We show that in the monitorable access to credit case the classical results are maintained. The contract will be renegotiation-proof, implementable via spot contracts and will coincide with the no-access optimum.

It is worth emphasizing that it has been noted in empirical work that contracts observed in reality have some characteristics that cannot be explained by the dynamic principal-agent models. In particular, payment schemes in which different outcomes imply equal reward are always theoretically suboptimal, although flat payment contracts are often observed in reality (Chiappori and Salanié (2000)). The systematical finding of downward wage constancy in empirical literature is also not explained either and the history dependence of multi-period contracts is not as strong as it is predicted by the classical theory (Bolton and Dewatripont , 2005). The predictions that our model presents as opposed to the classical models might be able to shed some light on some of this findings.

Finally, loss aversion induces a discontinuity in marginal utility. Therefore, the optimality conditions that are obtained in classical models are not valid since utility functions

are non-differentiable. We derive optimality conditions using convex analysis tools. In summary, the methodology is the following; we show that the program the Principal faces has a concave objective function and the feasible set is convex, therefore, the optimum can be characterized by a subdifferential equals zero condition (Rockafellar, 1974). The computation of the subdifferential in this case is not exempt of mathematical difficulty. It is derived inductively and, in order to so, we must obtain a new chain-rule that applies to the problem at hand. Note that several chain-rule formulas have been described in subdifferential calculus, but none of them applies to the functional forms of our model. Therefore, as a byproduct of this paper we introduce a new chain rule, which allow us to characterize the solutions to an optimization problem of this sort.

The following section presents the base model. Section 3 presents the optimality conditions in the full-commitment case and results that characterize the shape of an optimal payment scheme in our model. An analysis regarding the inter-temporal allocation of resources under the full-commitment is also carried out in Section 3, as well as a property regarding the shape of the optimal contract depending on the cost to the principal of the Participation Constraint. Section 4 focuses on the monitorable access to credit case. Section 5 presents a numerical example for two period optimal payment schemes, Graphical representations that illustrate how the optimal schemes vary with respect to parameters can be found in Section 5 as well. Section 6 presents the conclusions and discussions regarding our work.

## 1.2. The Base Model

The model is analogous to the dynamic moral hazard problem presented by Rogerson (1985) with the modification of the utility function to account for the reference dependent preferences. It consists in a repeated principal-agent problem with independent realizations of the outcome variable each period.

We assume that the relationship between the principal and the agent lasts $T$ periods. In each period the exerts an unobservable action $a_i \in \{a_L, a_H\}$ with $a_L < a_H$. The outcome in period $i$ is denoted $x_i \in [\underline{x}_i, \bar{x}_i]$ with a differentiable distribution function $f^i(x_i|a_i)$, where $a_i$ denotes the action chosen by the agent in period $i$. The distributions of outcomes are independent of each other. We assume that the distributions $f^i(x_i|a_i)$ have the MLRP property, that is, denoting $f^i_{a_i}(x_i|a_i) = f^i(x_i|a_H) - f^i(x_i|a_L)$, we must have $\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}$ is increasing in $x_i$. We let the wage schedule in period $i$ depend on the

outcomes obtained in all periods, denoted $\omega_i(x_0, x_1, \ldots, x_i)$ or $\omega_i(x_i)$ to simplify notation. The agent's utility function in period $i$ if she consumes $c_i$ and exerts the action $a_i$ is

$$\tilde{U}(c_i, R_i) - \psi_i(a_i) \tag{1.1}$$

$R_i$ is the reference point in period $i$, and $\psi_i(\cdot)$ is an increasing cost function. We denote $\Delta\psi_i = \psi_i(a_H) - \psi_i(a_L)$ We assume that $\tilde{U}$ is continuous.

Moreover, we assume the agent's preferences have the following property

$$\lim_{t \to 0^+} \frac{\tilde{U}_i(R+t, R) - \tilde{U}_i(R, R)}{t} < \lim_{t \to 0^+} \frac{\tilde{U}_i(R-t, R) - \tilde{U}_i(R, R)}{-t}$$

That is the left-sided derivative of $\tilde{U}_i$ is higher than its right-sided derivative. If the reference point is $R$ the marginal utility of payments over $R$ will be lower than the marginal loss in utility of an income that is below $R$ in the same amount. We assume that the agent's utility function is concave and that it is differentiable in all points other than $R$. In other words, the utility function presents loss aversion as introduced by Kahneman and Tversky's Prospect Theory (1979). In addition, since the utility function is non-differentiable at the reference point the classical analysis based on the first order conditions is no longer valid and in order to find optimality conditions one has to make use of convex analysis tools. In deriving the solution we use the concept of subdifferential and several of its properties as presented in Rockafellar (1974) and Bertsekas (2003).

We assume that the reference level depends dynamically on the consumption that took place in the previous period. The reference in each period corresponds to the consumption that took place in the previous period. This type of update is analogous to the one presented in Bowman et al. (1999) and Munro and Sugden (2003).[1]

Without loss of generality, following De Meza and Webb (2007) we assume the utility of period 0, $\tilde{U}_0$, takes the following form for $\ell_0 > 0$ and a smooth, concave and strictly increasing function $U(\cdot)$,

$$\tilde{U}_0(c_0, R_0) = U(c_0) - \ell_0\theta(c_0, R_0)(U(R_0) - U(c_0))$$

---

[1]Other papers assume a different formation process of the reference level. For instance, Köszegi and Rabin (2006) use the rational expectation of consumption. Gul (1991) takes the certainty equivalent as the reference. The issue of reference formation is still an understudied problem.

where

$$\theta(x, R) = \begin{cases} 1 & \text{if } x < R \\ 0 & \text{otherwise.} \end{cases} \tag{1.2}$$

The reference in period 0, $R_0$ is exogenously given. A graphical representation of the first period utility fuction is given in Figure 1.2. As the reference increases from $R'$ to $R$ the utility above $R$ remains unaffected, while it drops below $R$.[2]

In the classical models of Dynamic Moral Hazard the agent has the same utility function in all periods. In our model, the utility of each period depends on the reference level of the agent, and at the same time, retains its functional form in order to represent the utility of the same agent across periods. We assume that the functional form of the utility for consumption levels above the reference remains the same from one period to the next. For consumptions under the reference the utility is decreasing in the reference. The reference level in period $i + 1$ corresponds to $c_i$. The utility in period $i + 1$ takes the following form,

$$\tilde{U}_{i+1}(c_{i+1}, R_i) = \tilde{U}_{i+1}(c_{i+1}, c_i) = U(c_{i+1}) - \ell_{i+1}\theta(c_{i+1}, c_i)\left(U(c_i) - U(c_{i+1})\right)$$

with $\theta(c_{i+1}, c_i)$ defined as (1.2) and $\ell_{i+1} > 0$. We assume that $\ell_i\delta < 1$ in order to assure that the total utility of two consecutive periods is increasing in the consumption of the first of the two periods. [3]

The principal is risk neutral and therefore for a given outcome $x_i$ her utility will be $x_i - \omega_i(x_i)$. The agent and the principal discount the future at factor $\delta$. Also, we make the standard assumption that unlimited transfers of utility from the principal to the agent are possible in every period. This last assumption is necessary to prove a result in the monitorable access to credit section.

---

[2]Any function $\tilde{U}$ that presents a kink can be rewritten, for some differentiable, increasing and concave function $U$ and positive constant $\ell$, as

$$\tilde{U}(c, R) = U(c) - \ell\theta(c, R)\left(U(R) - U(c)\right)$$

$\theta(c, R) = 1$ if $c \leq R$ and $\theta(c, R) = 0$ otherwise.

[3]For generality we let $\ell_i$ vary across periods. For instance, in Harrison and List (2004) it is noted that loss aversion is mitigated by market experience, suggesting that $\ell_i$ might decrease over time.

Figure 1.1: Utility function for different levels of reference.



## 1.3. Dynamic contracts with no access to credit markets

In the first part of this section we find optimality conditions for the dynamic moral hazard problem under full-commitment and no-access to credit. In order to do this we define a utility provision cost function $h_i(v_0, v_1, \ldots, v_i)$ that represents the cost to the principal of providing a value of utility $v_i$ in period $i$ if the utility provisions in the previous periods were $\{v_0, v_1, \ldots, v_{i-1}\}$. Because of the reference dependent preferences and dynamic update assumptions, the utility provision cost function in each period depends on the utility provisions of all previous periods. We show that $h_i$ is convex in $\mathbb{R}^i$ which implies the concavity of the objective function of the principal. Next we

find the subdifferential of the function $h_i$, by writing it as a composition of functions. A chain rule for this particular case is derived to compute the subdifferential of the compositions. Finally, a general formula for the subdifferential is obtained inductively. Using the computations described and results of subdifferential calculus, we obtain the subdifferential of the Lagrangian pertaining to the optimization problem that the principal faces. The optimality conditions are thus derived by making the latter equal to zero.

The second part of this section describes some characteristics of the optimal payment scheme that are implied by the optimality conditions. The payment scheme has a flat segment in every period after the first one. This flat segment may extend for the entire support of outcomes in periods before the last one. In the last period, for every set of outcomes in previous periods, the payment scheme must be contingent on the results of all periods.

Next we analyze the inter-temporal properties of the optimal payment schemes. We emphasize that in addition to the fact that payments schemes are not be strictly increasing with respect to $x_i$ in any period $i > 0$, they may not be strictly increasing (as functions) in periods after $i$ with respect to period $i$'s outcome $x_i$. This suggests that our model may predict a smaller dependence on payments schemes on outcomes across and within periods in comparison with the canonical model. We find a relationship between a period's payment and future period's payment. We explain why it differs to the analogous relationship from the canonical model in Rogerson (1985).

We next analyze whether in any period $i$ the agent would like to borrow or save for the next period. We note first that the preferences present a "status quo bias". Situations may arise that the agent would lose utility either by saving or borrowing and the agent is compelled to consume his allocation. Furthermore, for any set of payments in two consecutive periods the smallest interest rate at which he is willing to lend part of his income is strictly higher than the interest rate at which he is willing to borrow. This is not the case in the canonical model since, indifference between lending and borrowing, for any set of payments schemes for two consecutive periods, can be attained for a single interest rate. This gap between willingness to accept and willingness to pay has been described in other economical contexts (Bateman et al. (1997)).

We then determine whether the agent would like to borrow or save for some possi-

ble shapes of payment schemes that can be obtained in our model. We find that if the payment scheme is such that the agent is paid the reference for low outcomes and is rewarded for good outcomes the agent will not have incentives to save.

Finally, we show that as in the canonical model the optimal payment scheme is renegotiation-proof and is not spot implementable. The optimal sequence of spot contracts exhibits memory and does not coincide with the full-commitment optimum in general.

### 1.3.1. Full commitment case

If there is no access to credit markets the agent's consumption in period each period will equal the payment the agent receives. Therefore, the principal faces the following program,

$$\max_{(\omega_i(\cdot))_i,(a_i)_i} \sum_{i=0}^{T} \delta^i \mathbb{E}\left(x_i - \omega_i(x_0, x_1, \dots, x_i)|a_0, a_1, \dots, a_i\right)$$

subject to

$$\sum_{i=0}^{T} \delta^i \left(\mathbb{E}\left(\tilde{U}_i(\omega_i(x_0, x_1, \dots, x_i), c_{i-1})|a_0, a_1, \dots, a_i\right) - \psi_i(a_i)\right) \geq U^* \quad (PC)$$

$$a = (a_0, a_1(x_1), \dots a_T(x_0, x_1, \dots, x_T)) \in \text{argmax}_a \sum_{i=0}^{T} \delta^i \left(\mathbb{E}\left(\tilde{U}_i(\omega_i(x_0, x_1, \dots, x_i), c_{i-1})|a_0, a_1, \dots, a_i\right) - \psi(a_i)\right) \quad (IC)$$

Where $\mathbb{E}(\cdot|a_0, a_1, \dots, a_i)$ is the expectation of given that the agent chooses the actions $(a_0, a_1, \dots, a_i)$, and for any function $g$ is given by,

$$\mathbb{E}\left(g(x_0, x_1, \dots, x_i)|a_0, a_1, \dots, a_i\right) = \int \int \cdots \int g(x_0, x_1, \dots, x_i) f^1(x_0|a_0) f^2(x_1|a_1) \cdots f^i(x_i|a_i) dx_0 dx_1 \dots dx_i.$$

The objective function represents the expected payment the principal will get from the contract, the first constraint (PC) is the participation constraint (PC) and it requires that the agent gets an expected utility of a at least $U^*$ from the relationship with the principal. The constraint (IC) states that the effort chosen maximizes the expected utility of the agent, and is henceforth referred to as the incentive compatibility constraint (IC).

If the Principal wants to implement the high effort in each period the last constraint will be,

$$\sum_{j=i}^{T} \delta^j (\mathbb{E}\left(\tilde{U}_j(\omega_j(x_0, x_1, \dots, x_j), c_{j-1})|a_i = a_H, \dots, a_j\right) -$$

$$\mathbb{E}\left(\tilde{U}_j(\omega_j(x_0, x_1, \dots, x_j), c_{j-1})|a_i = a_L, \dots, a_j\right)) - \Delta\psi_i = 0 \quad \forall(x_0, x_1, \dots, x_{i-1}) \quad (1.3)$$

18

that is

$$\sum_{j=i}^{T} \delta^j \int \tilde{U}_j(\omega_j(x_0, x_1, \ldots, x_j), c_{j-1}) f_{a_i}^i(x_i|a_i) \cdots f^j(x_j|a_j) dx_0 dx_1 \ldots dx_j - \Delta \psi_i = 0 \quad \forall (x_0, x_1, \ldots, x_{i-1}) \qquad \text{(ICi)}$$

### 1.3.1.1. Characterization of the optimal payment scheme

We find that for all periods there may be an interval of outcomes for which the optimal payment corresponds to the period's reference. Since the utility function presents loss aversion, for an outcome in which the optimal payment is at the reference it can be costly for the Principal to give a payment strictly over or under the reference for outcomes slightly better or worse. A payment over the reference will provide a relatively low marginal utility and therefore create low incentives. A payment under the reference will decrease the agent's utility quickly, providing incentives, but straining the PC. Therefore, a payment scheme that gives the reference for an outcome, pays the reference for close outcomes as well. Consumption smoothing requires that the reference be reached from period 1 on. In those periods, the optimal payment scheme must indeed have a flat segment.

The utility function is non-differentiable and therefore we must derive the optimal payment scheme using tools from convex analysis. In particular, we use the concept of subdifferential, which is an extension of the common concept of the differential that is commonly used to solve economic problems. In order to find optimality conditions, it is convenient to rewrite the program the principal faces as follows,

$$\max_{(v_i(\cdot))_i, (a_i)_i} \sum_{i=0}^{T} \delta^i \mathbb{E} \left( x_i - h_i(v_0(x_0), v_1(x_0, x_1), \ldots, v_i(x_0, x_1, \ldots, x_i)) | a_0, a_1, \ldots, a_i \right) \qquad (1.4)$$

subject to

$$\sum_{i=0}^{T} \delta^i \left( \mathbb{E} \left( v_i(x_0, x_1, \ldots, x_i) | a_0, a_1, \ldots, a_i \right) - \psi_i(a_i) \right) \geq U^* \qquad \text{(PC')}$$

$$a = (a_0, a_1(x_1), \ldots a_T(x_0, x_1, \ldots, x_T)) \in \text{argmax}_a \sum_{i=0}^{T} \delta^i \left( \mathbb{E} \left( v_i(x_0, x_1, \ldots, x_i) | a_0, a_1, \ldots, a_i \right) - \psi(a_i) \right)$$

$$\text{(IC')}$$

where $h_i(v_0, v_1, \ldots, v_i)$ represents the cost of providing utility $v_i$ in period $i$. This cost depends on the previous utility provisions because of the reference dependence of the agent's utility and is an increasing and continuous function given by,[4]

---

[4]$h_i$ is obtained inverting the utility provision $v_i = U(h_i) - \ell_i \theta(h_i, h_{i-1}) \left( U(h_{i-1}) - U(h_i) \right)$ with respect to the payment $h_i$.

$$h_i(v_0, v_1, \ldots, v_i) = \begin{cases} U^{-1}(v_i) & \text{if } v_i \geq U(h_{i-1}(v_0, v_1, \ldots, v_{i-1})) \\ U^{-1}\left(\frac{v_i + \ell_i U(h_{i-1}(v_0, v_1, \ldots, v_{i-1}))}{1+\ell_i}\right) & \text{if } v_i < U(h_{i-1}(v_0, v_1, \ldots, v_{i-1})) \end{cases} \quad (1.5)$$

where $U(h_{-1}) = R_0$.

**Convexity, subdifferential and optimality conditions**

**Property 1** (Convexity). *Under the assumptions of the model the utility provision cost functions $h_i : \mathbb{R}^i \rightarrow \mathbb{R}$ for $i \in 1, \ldots, T$ are strictly convex and therefore the optimization problem given by (1.4)-(PC')-(IC') has a strictly concave objective function and the feasible set is convex.*

*Proof.* Let's see that $h_i(v_0, v_1, \ldots, v_i)$ is strictly convex. In order to do that, let's first note that $h_i(v_0, v_1, \ldots, v_i) = U^{-1}(U(h_i(v_0, v_1, \ldots, v_i)))$, we prove that $U(h_i(v_0, v_1, \ldots, v_i))$ is strictly convex and increasing and we conclude by the strict convexity of $U^{-1}$ (implied by the strict concavity of $U$).[5] Let $(v_0, \ldots, v_i)$ and $(v'_0, \ldots, v'_i)$ be two utility provision vectors, we prove that

$$U(h_i(\lambda(v_0, \ldots, v_i) + (1-\lambda)(v'_0, \ldots, v'_i))) < \lambda U(h_i(v_0, v_1, \ldots, v_i)) + (1-\lambda)U(h_i(v'_0, v'_1, \ldots, v'_i)) \quad \forall \lambda \in (0,1) \quad (1.6)$$

Note that for $i = 0$, by (1.5) $U(h(v_0))$ is linear by parts, increasing and convex. Let's prove (1.6) assuming true for $i - 1$.

If $\lambda v_i + (1 - \lambda)v'_i < U(h_{i-1}(\lambda(v_0, \ldots, v_{i-1}) + (1 - \lambda)(v'_0, \ldots, v'_{i-1})))$ then

$$
\begin{aligned}
U(h_i(\lambda(v_0, \ldots, v_i) + (1-\lambda)(v'_0, \ldots, v'_i))) &= \left(\frac{\lambda v_i + (1-\lambda)v'_i + \ell_i U(h_{i-1}(\lambda(v_0, \ldots, v_{i-1}) + (1-\lambda)(v'_0, \ldots, v'_{i-1})))}{1+\ell_i}\right) \\
&\leq \lambda\left(\frac{v_i + \ell_i U(h_{i-1}(v_0, v_1, \ldots, v_{i-1}))}{1+\ell_i}\right) + \\
&\quad (1-\lambda)\left(\frac{v'_i + \ell_i U(h_{i-1}(v'_0, v'_1, \ldots, v'_{i-1}))}{1+\ell_i}\right) \\
&\leq \lambda U(h_i(v_0, v_1, \ldots, v_i)) + (1-\lambda)U(h_i(v'_0, v'_1, \ldots, v'_i))
\end{aligned}
$$

Where the first inequality is implied by the induction hypothesis and the second is justified noting that if $v_i > U(h_{i-1}(v_0, v_1, \ldots, v_{i-1}))$ then $v_i > \left(\frac{v_i + \ell_i U(h_{i-1}(v_0, v_1, \ldots, v_{i-1}))}{1+\ell_i}\right)$ and if $v_i \leq U(h_{i-1}(v_0, v_1, \ldots, v_{i-1}))$ then $v_i = \left(\frac{v_i + \ell_i U(h_{i-1}(v_0, v_1, \ldots, v_{i-1}))}{1+\ell_i}\right)$.

A similar argument proves (1.6) for the case in which $\lambda v_i + (1-\lambda)v'_i \geq U(h_{i-1}(\lambda(v_0, \ldots, v_{i-1}) +$

---

[5]Note that the composition of a convex increasing function with a convex function is convex.

$(1-\lambda)(v'_0, \ldots, v'_{i-1})))$. Finally it is straightforward to verify the convexity of the feasible set.

$\square$

**Property 2** (subdifferential). *The subdifferential[6] of $h_i(v_0, v_1, \ldots, v_i)$ is given by*

$$\partial h_i(v_0, v_1, \ldots, v_i) = \left( \frac{1}{U'(\omega_i)} \left( \prod_{t=j+1}^{i} \frac{k_t(x_0, x_1, \ldots, x_t)\ell_t}{1 + k_t(x_0, x_1, \ldots, x_t)\ell_t} \right) \frac{1}{1 + k_j(x_0, x_1, \ldots, x_j)\ell_j} \right)_{j=0}^{i} \tag{1.7}$$

*where*

$$k_t(x_0, x_1, \ldots, x_t) \in \begin{cases} \{1\} & \text{if } \omega_t(x_0, x_1, \ldots, x_t) < R_t \\ [0,1] & \text{if } \omega_t(x_0, x_1, \ldots, x_t) = R_t \\ \{0\} & \text{otherwise} \end{cases} \tag{1.8}$$

*Proof.* We have

$$h_i(v_0, v_1, \ldots, v_i) = U^{-1}(U(h_i(v_0, v_1, \ldots, v_i)))$$

where

$$U(h_i(v_0, v_1, \ldots, v_i)) = \begin{cases} v_i & \text{if } v_i \geq U(h_{i-1}(v_0, v_1, \ldots, v_{i-1})) \\ \frac{v_i + \ell_i U(h_{i-1}(v_0, v_1, \ldots, v_{i-1}))}{1 + \ell_i} & \text{if } v_i < U(h_{i-1}(v_0, v_1, \ldots, v_{i-1})) \end{cases} \tag{1.9}$$

By Proposition 4.2.5 in Bertsekas (2003) we know that

$$\partial h_i(v_0, v_1, \ldots, v_i) = \left( U^{-1} \right)' \left( (U \circ h_i)(v_0, v_1, \ldots, v_i) \right) \cdot \partial \left( (U \circ h_i)(v_0, v_1, \ldots, v_i) \right)$$

.

Now, note that from (1.9) we have $U(h_i(v_0, v_1, \ldots, v_i)) = F_i((U \circ h_{i-1})(v_0, v_1, \ldots, v_{i-1}), v_i)$ where $F_i(x, y) = \begin{cases} y & \text{if } y \geq x \\ \frac{y + \ell_i x}{1 + \ell_i} & \text{if } y < x \end{cases}$

Let $(d_0, \ldots, d_{i-1}) \in \partial (U \circ h_{i-1})(v_0, v_1, \ldots, v_{i-1})$ and $(\tilde{d}_0, \tilde{d}_1) \in \partial F_i((U \circ h_{i-1})(v_0, v_1, \ldots, v_{i-1}), v_i)$ let's see that that $(d_0 \cdot \tilde{d}_0, d_1 \cdot \tilde{d}_0, \ldots, d_{i-1} \cdot \tilde{d}_0, \tilde{d}_1) \in \partial (U \circ h_i)(v_0, v_1, \ldots, v_i)$. In fact, we

---

[6]By definition we know that for a generic convex function $f : \mathbb{R}^{n+1} \to \mathbb{R}$ the subdifferential at $x \in \mathbb{R}^{n+1}$ will be given by the set of vectors $d = (d_0, d_1, \ldots, d_{n+1}) \in \mathbb{R}^{n+1}$ such that for any vector $\alpha \in \mathbb{R}^{n+1}$

$$f(x + \alpha) \geq f(x) + d \cdot \alpha$$

have

$$(U \circ h_i)(v_0 + \alpha_0, v_1 + \alpha_1, \ldots, v_i + \alpha_i) =$$

$$F_i((U \circ h_{i-1})(v_0 + \alpha_0, \ldots, v_{i-1} + \alpha_{i-1}), v_i + \alpha_i) \geq F_i((U \circ h_{i-1})(v_0, v_1, \ldots, v_{i-1}) + d_0\alpha_0 + \cdots + d_{i-1}\alpha_{i-1}, v_i + \alpha_i)$$

$$\geq F_i((U \circ h_{i-1})(v_0, v_1, \ldots, v_{i-1}), v_i) + d_0\tilde{d}_0\alpha_0 + \cdots + d_{i-1}\tilde{d}_0\alpha_{i-1} + \tilde{d}_1\alpha_i$$

Where the first inequality is due to $(d_0, \ldots, d_{i-1}) \in \partial(U \circ h_{i-1})(v_0, v_1, \ldots, v_{i-1})$ and $F_i$ increasing in its first variable. The second inequality is implied by

$$(\tilde{d}_0, \tilde{d}_1) \in \partial F_i((U \circ h_{i-1})(v_0, v_1, \ldots, v_{i-1}), v_i)$$

. We see now that the reverse is also true.

Let's compute $\partial F_i(x, y)$. In the points $x \neq y$ $F_i$ is differentiable and therefore its subdifferential coincides with the derivative.

Otherwise, $y = x$ and the elements of the subdifferential of $\partial F_i(x, y)$ will be the pairs $(\tilde{d}_0, \tilde{d}_1)$ such that,

$$F_i(x + \alpha_0, y + \alpha_1) \geq F(x, y) + \alpha_0\tilde{d}_0 + \alpha_1\tilde{d}_1 \qquad \forall \alpha_0, \alpha_1 \in \mathbb{R} \tag{1.10}$$

If $\alpha_0 \leq \alpha_1$ then $x + \alpha_0 \leq y + \alpha_1$ and (1.10) becomes $(1 - \tilde{d}_1)\alpha_1 \geq \tilde{d}_0\alpha_0$ which is true for all $\alpha_1 \geq \alpha_0$ if and only if $(1 - \tilde{d}_1) = \tilde{d}_0 > 0$.

If $\alpha_0 > \alpha_1$ then $x + \alpha_0 > y + \alpha_1$ and (1.10) becomes $\left(\frac{\ell_i}{1+\ell_i} - \tilde{d}_0\right)\alpha_0 \geq \left(\tilde{d}_1 - \frac{1}{1+\ell_i}\right)\alpha_1$ which is true for all $\alpha_0 < \alpha_1$ if and only if $\left(\frac{\ell_i}{1+\ell_i} - \tilde{d}_0\right) = \left(\tilde{d}_1 - \frac{1}{1+\ell_i}\right) > 0$. Therefore, $\tilde{d}_0 \in \left[0, \frac{\ell_i}{1+\ell_i}\right]$, $\tilde{d}_1 \in \left[\frac{1}{1+\ell_i}, 1\right]$ and $\tilde{d}_0 = 1 - \tilde{d}_1$.

We conclude that,

$$\partial F_i(x, y) = \left\{ \left(\frac{k_i\ell_i}{1 + \ell_i k_i}, \frac{1}{1 + \ell_i k_i}\right); \quad \text{where} \quad k_i(x_0, x_1, \ldots, x_i) \in \left\{ \begin{array}{ll} \{1\} & \text{if } \omega_i(x_0, x_1, \ldots, x_i) < R_i \\ [0,1] & \text{if } \omega_i(x_0, x_1, \ldots, x_i) = R_i \\ \{0\} & \text{otherwise} \end{array} \right. \right\}$$

Suppose that $(\bar{d}_0, \ldots, \bar{d}_i) \in \partial(U \circ h_i)(v_0, \ldots, v_i)$, let's see that $(\bar{d}_0, \bar{d}_1, \ldots \bar{d}_{i-1}, \bar{d}_i) = (d_0 \cdot \tilde{d}_0, d_1 \cdot \tilde{d}_0, \ldots, d_{i-1} \cdot \tilde{d}_0, \tilde{d}_1)$ for some vectors $(d_0, \ldots, d_{i-1}) \in \partial(U \circ h_{i-1})(v_0, \ldots, v_{i-1})$ and $(\tilde{d}_0, \tilde{d}_1) \in \partial F_i((U \circ h_{i-1})(v_0, \ldots, v_{i-1}), v_i)$.

From $F_i((U \circ h_{i-1})(v_0, \ldots, v_{i-1}), v_i + \alpha_i) \geq F_i((U \circ h_{i-1})(v_0, \ldots, v_{i-1}), v_i) + \alpha_i\bar{d}_i \quad \forall \alpha_i$ we must have $\bar{d}_i = \frac{1}{1+\ell_i k_i}$ with $k_i$ defined by (1.8) (subdifferential in one variable). We know that in points in which $F_i$ is differentiable its subdifferential coincides with the derivative

which will be $(0,1)$ if $(U \circ h_{i-1})(v_0, \ldots, v_{i-1}) < v_i$ and $\left(\frac{\ell_i}{1+\ell_i}, \frac{1}{1+\ell_i}\right)$ if $(U \circ h_{i-1})(v_0, \ldots, v_{i-1}) > v_i$. Therefore, using Proposition 4.2.5 Bertsekas (2003) we must have that $\partial U(h_i(v_0, \ldots, v_i)) = \left((1-\bar{d}_i) \cdot \partial(U \circ h_{i-1})(v_0, \ldots, v_{i-1}), \bar{d}_i\right)$.

If $F$ is not differentiable we have $(U \circ h_{i-1})(v_0, \ldots, v_{i-1}) = v_i$. Let $(\alpha_0, \alpha_1, \ldots, \alpha_{i-1}) \in \mathbb{R}^i$, we define $\hat{\alpha}_i = (U \circ h_{i-1})(v_0 + \alpha_0, \ldots, v_{i-1} + \alpha_{i-1}) - v_i$. Since $(\bar{d}_0, \ldots, \bar{d}_i)$ in $\partial U(h_i(v_0, \ldots, v_i))$ we have

$$v_i + \hat{\alpha}_i \geq v_i + \bar{d}_0\alpha_0 + \cdots + \bar{d}_{i-1}\alpha_{i-1} + \bar{d}_i\hat{\alpha}_i$$

$$\implies \hat{\alpha}_i(1-\bar{d}_i) \geq \bar{d}_0\alpha_0 + \cdots + \bar{d}\alpha_{i-1}$$

$$\implies (U \circ h_{i-1})(v_0 + \alpha_0, \ldots, v_{i-1} + \alpha_{i-1}) - (U \circ h_{i-1})(v_0, \ldots, v_{i-1}) \geq \left(\bar{d}_0\alpha_0 + \cdots + \bar{d}_{i-1}\alpha_{i-1}\right)\frac{1}{(1-\bar{d}_i)}$$

$$\implies (\bar{d}_0, \ldots, \bar{d_{i-1}})\frac{1}{(1-\bar{d}_i)} \in \partial(U \circ h_{i-1})(v_0, \ldots, v_{i-1})$$

We conclude for $(d_0, \ldots, d_{i-1}) = (\bar{d}_0, \ldots, \bar{d}_i)\frac{1}{(1-d_i)} \in \partial(U \circ h_{i-1})(v_0, \ldots, v_{i-1})$ and $(\tilde{d}_0, \tilde{d}_1) = \left((1-\bar{d}_i), \bar{d}_i\right) \in \partial F((U \circ h_{i-1})(v_0, \ldots, v_{i-1}), v_i)$ that $(\bar{d}_0, \bar{d}_1, \ldots \bar{d}_{i-1}, \bar{d}_i) = (d_0 \cdot \tilde{d}_0, d_1 \cdot \tilde{d}_0, \ldots, d_{i-1} \cdot \tilde{d}_0, \tilde{d}_1)$.

We deduce inductively $\partial U(h_i(v_0, \ldots, v_i))$. For the functions $k_i$ defined by (1.8) we have

$$\partial U(h(v_0)) = \left\{\frac{1}{1+k_0\ell_0}\right\}$$

therefore

$$\partial U(h(v_0, v_1)) = \left(\frac{k_1\ell_1}{1+\ell_1 k_1} \cdot \frac{1}{1+k_0\ell_0}, \frac{1}{1+\ell_1 k_1}\right)$$

and

$$\partial U(h(v_0, v_1, v_2)) = \left(\frac{k_2\ell_2}{1+\ell_2 k_2} \cdot \frac{k_1\ell_1}{1+\ell_1 k_1} \cdot \frac{1}{1+k_0\ell_0}, \frac{k_2\ell_2}{1+\ell_2 k_2} \cdot \frac{1}{1+k_1\ell_1}, \frac{1}{1+\ell_2 k_2}\right)$$

and inductively, (1.7) is obtained.

$\square$

**Property 3** (Optimality Conditions). *There is a unique optimal wage schedule that solves the program faced by the principal and it is characterized by the following optimality conditions,*

$$\frac{1}{U'(\omega_i(x_0, x_1, \ldots, x_i))} = (1+k_i(x_0, x_1, \ldots, x_i)\ell_i)\left(\lambda_i + \mu_i \frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}\right) +$$

$$-\delta\ell_{i+1}\int_{\omega_{i+1}\leq\omega_i} k_{i+1}(x_0, x_1, \ldots, x_{i+1})(\lambda_{i+1}+\mu_{i+1}\frac{f^{i+1}_{a_{i+1}}(x_{i+1}|a_{i+1})}{f^{i+1}(x_{i+1}|a_{i+1})})f^{i+1}(x_{i+1}|a_{i+1})dx_{i+1}. \qquad \forall i < T$$

$$(1.11)$$

*and*

$$\frac{1}{U'(\omega_T(x_0, x_1, \ldots, x_T))} = (1 + k_T(x_0, x_1, \ldots, x_T)\ell_T)\left(\lambda_T + \mu_T \frac{f_{a_T}^T(x_T|a_i)}{f^i(x_i|a_i)}\right) \qquad (1.12)$$

*where*

- $\lambda_i = \lambda + \sum_{k=0}^{i-1} \mu_k \frac{f_{e_k}^k(x_k|a_k)}{f^k(x_k|a_k)}$, *with* $\lambda$ *the multiplier associated to (PC') and* $\mu_i = \mu_i(x_0, \ldots, x_{i-1})$ *the multiplier associated to 1.3.*

- *The function* $k_i(x_0, x_1, \ldots, x_i)$ *is associated to the kink in the utility function and is given by,*

$$k_i(x_0, x_1, \ldots, x_i) \in \begin{cases} \{1\} & \text{if } \omega_i(x_0, x_1, \ldots, x_i) < R_i \\ [0,1] & \text{if } \omega_i(x_0, x_1, \ldots, x_i) = R_i \\ \{0\} & \text{otherwise} \end{cases}$$

*Proof.* We assume that the principal is looking for optimal utility provisions $(v_i(\cdot))_{i=0}^T$ in the function spaces $\left(L^1([\underline{x}_0, \bar{x}_0] \times [\underline{x}_1, \bar{x}_1] \times \cdots \times [\underline{x}_i, \bar{x}_i])\right)_{i=0}^T$. By Karush-Kuhn-Tucker conditions the optimal payment scheme will be the maximum of the the Lagrangian function L given by[7]

$$L = \sum_{i=0}^T \delta^i \mathbb{E}\left(x_i - h_i(v_0(x_o), v_1(x_0, x_1), \ldots, v_i(x_0, x_1, \ldots, x_i))|a_0, a_1, \ldots, a_i\right) +$$

$$\lambda\left(\sum_{i=0}^T \delta^i \left(\mathbb{E}\left(v_i(x_0, x_1, \ldots, x_i)|a_0, a_1, \ldots, a_i\right) - \psi_i(a_i)\right) - U^*\right) +$$

$$\sum_{i=0}^T \left(\sum_{j=i}^T \delta^j \left(\Delta_{a_i} \mathbb{E}\left(\mu_i(x_0, x_1, \ldots, x_i) \cdot v_j(x_0, x_1, \ldots, x_j)|a_0, a_1, \ldots, a_j\right)\right) - \Delta\psi_i\right)$$

Where we denote

$$\Delta_{a_i} \mathbb{E}\left(\mu_i(x_0, x_1, \ldots, x_i) \cdot v_j(x_0, x_1, \ldots, x_j)|a_i, \ldots, a_j\right) =$$

$$\mathbb{E}\left(\mu_i(x_0, x_1, \ldots, x_i) \cdot v_j(x_0, x_1, \ldots, x_j)|a_i = a_H, \ldots, a_j\right) +$$

$$-\mathbb{E}\left(\mu_i(x_0, x_1, \ldots, x_i) \cdot v_j(x_0, x_1, \ldots, x_j)|a_i = a_L, \ldots, a_j\right) =$$

$$\int \mu_i(x_0, x_1, \ldots, x_i) \cdot v_j(x_0, x_1, \ldots, x_j) f_{a_i}^i(x_i|a_i) \cdots f^j(x_j|a_j) dx_0 dx_1 \ldots dx_j$$

$L$, although non-differentiable, is concave in $(v_i(\cdot))_i$ and the set of constraints is convex (Property 1), therefore a necessary and sufficient condition for a wage schedule to be optimal is that it makes a subdifferential of $-L$ (denoted $\partial(-L)$) equal to 0. Since from Property 1 we have that $\mathbb{E}\left(h_i(v_0(x_o), v_1(x_0, x_1), \ldots, v_i(x_0, x_1, \ldots, x_i))|a_0, a_1, \ldots, a_i\right)$

---

[7]Note that there is an infinite number of contraints since 1.3 must be fullfilled $\forall(x_0, x_1, \ldots, x_i)$. See Rockafellar (1974) for details on how to derive a Lagrangian in this case.

are convex in $(v_i(\cdot))_i$, from proposition 4.2.4 in Bertsekas (2003). [8]

$$\partial(-L) = \sum_{i=0}^{T} \delta^i \partial \mathbb{E}\left(h_i(v_0(x_o), v_1(x_0, x_1), \ldots, v_i(x_0, x_1, \ldots, x_i))|a_0, a_1, \ldots, a_i\right) +$$

$$- \lambda \left(\sum_{i=0}^{T} \delta^i \partial \mathbb{E}\left(v_i(x_0, x_1, \ldots, x_i)|a_0, a_1, \ldots, a_i\right)\right) +$$

$$- \sum_{i=0}^{T} \left(\sum_{j=i}^{T} \delta^j \partial \left(\Delta_{a_i} \mathbb{E}\left(\mu_i(x_0, x_1, \ldots, x_i) \cdot v_j(x_0, x_1, \ldots, x_j)|a_0, a_1, \ldots, a_j\right)\right)\right)$$

From Theorem 22 of Rockafellar (1974), we know that a subdifferential of $-L$ is the expectation of the subdifferential of the integrand.

$$\partial(-L) = \sum_{i=0}^{T} \delta^i \mathbb{E}\left(\partial h_i(v_0(x_o), v_1(x_0, x_1), \ldots, v_i(x_0, x_1, \ldots, x_i))|a_0, a_1, \ldots, a_i\right) +$$

$$- \lambda \left(\sum_{i=0}^{T} \delta^i \mathbb{E}\left(\partial v_i(x_0, x_1, \ldots, x_i)|a_0, a_1, \ldots, a_i\right)\right) +$$

$$- \sum_{i=0}^{T} \left(\sum_{j=i}^{T} \delta^j \left(\Delta_{a_i} \mathbb{E}\left(\mu_i(x_0, x_1, \ldots, x_i)|\partial v_j(x_0, x_1, \ldots, x_j)|a_0, a_1, \ldots, a_j\right)\right)\right)$$

Therefore from Property 2 making $\partial(-L)$ equal 0 by components corresponds to,

$$0 = \sum_{i=j}^{T} \delta^{i-j} \mathbb{E}\left(\frac{1}{U'(\omega_i(x_0, x_1, \ldots, x_i))}\left(\prod_{t=j+1}^{i} \frac{k_t(x_0, x_1, \ldots, x_t)\ell_t}{1 + k_t(x_0, x_1, \ldots, x_t)\ell_t}\right) \frac{1}{1 + k_j(x_0, x_1, \ldots, x_j)\ell_j} \cdot g(x_0, x_1, \ldots, x_j)|a_0, a_1, \ldots, a_i\right) +$$

$$- \lambda \mathbb{E}\left(g(x_0, x_1, \ldots, x_j)|a_0, a_1, \ldots, a_j\right) +$$

$$- \sum_{i=0}^{j} \left(\Delta_{a_i} \mathbb{E}\left(\mu_i(x_0, x_1, \ldots, x_i)g(x_0, x_1, \ldots, x_j)|a_0, a_1, \ldots, a_j\right)\right)$$

for every $g \in L^1\left([\underline{x_0}, \bar{x_0}] \times [\underline{x_1}, \bar{x_1}] \times \cdots \times [\underline{x_j}, \bar{x_j}]\right)$, which implies

$$\frac{1}{U'(\omega_j)} \cdot \frac{1}{1 + k_j \ell_j} + \sum_{i=j+1}^{T} \delta^{i-j} \mathbb{E}\left(\frac{1}{U'(\omega_i)}\left(\prod_{t=j+1}^{i} \frac{k_t \ell_t}{1 + k_t \ell_t}\right) \frac{1}{1 + k_j \ell_j}|a_{j+1}, \ldots, a_i\right) = \lambda + \sum_{i=0}^{j} \mu_i \frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)} \quad (1.13)$$

$$= \lambda_j + \mu_j \frac{f^j_{a_j}(x_j|a_j)}{f^j(x_j|a_j)}$$

For every $j \in \{1, \ldots, T\}$. Multiplying the equation for $j+1$ by $k_{j+1}\ell_{+1}$ taking expectation with respect to $f^{j+1}(x_{i+1}|a_{j+1})$ we obtain

---

[8]Note that the subdifferential of a constant equals 0

$$\mathbb{E}\left(\frac{1}{U'(\omega_{j+1})}\frac{k_{j+1}\ell_{j+1}}{1+k_{j+1}\ell_{j+1}}|a_{j+1}\right) = -\sum_{i=j+2}^{T}\delta^{i-j-1}\mathbb{E}\left(\frac{1}{U'(\omega_i)}\left(\prod_{t=j+2}^{i}\frac{k_t\ell_t}{1+k_t\ell_t}\right)\frac{k_{j+1}\ell_{j+1}}{1+k_{j+1}\ell_{j+1}}|a_{j+1},\ldots,a_i\right) +$$

$$\mathbb{E}\left(\left(\lambda_{j+1}+\mu_{j+1}\frac{f_{a_{j+1}}^{j+1}(x_{j+1}|a_{j+1})}{f^{j+1}(x_{j+1}|a_{j+1})}\right)k_{j+1}\ell_{j+1}|a_{j+1}\right)$$

Replacing this last expression in the $j+1$ term of the sum in (1.13), 1.11 and 1.21 are obtained.

$\square$

Making $\ell_i = 0 \quad \forall i$ we obtain the optimality condition with a differentiable utility function, which we refer to as the "classical case". The optimality condition for a spot contract will be given by (1.21), with $T = 0$, as obtained by De Meza and Webb (2007).

**Shape of the optimal payment scheme**

The following property gives a further characterization of the optimal payment schemes in our model,

**Property 4** (Shape of the optimal payment scheme). *If $\ell_{i-1} \geq \ell_i \geq \ell_{i+1}$ and $\delta\ell_i \leq 1$, then $\omega_i(x_0, x_1, \ldots, x_i)$ is continuous and non-decreasing in $x_i$ and,*

1. *For $i \in \{1, \ldots, T\}$, if $\omega_{i-1}(x_0, x_1, \ldots, x_{i-1}) > R_{i-1}$ then for any value of $(x_0, x_1, \ldots, x_{i-1})$ we must have $\omega_i(x_0, x_1, \ldots, x_i) = R_i$ for some outcome $x_i \in [\underline{x}_i, \bar{x}^i]$ . Furthermore the payment scheme has a flat segment at the reference and therefore, $\omega_i$ is not be strictly increasing.*

2. *For $i \in \{1, \ldots, T-1\}$, if $(\ell_i - \ell_{i+1})\delta \geq \ell_{i-1} - \ell_i$ and $\omega_{i-1}(x_0, x_1, \ldots, x_{i-1}) \leq R_{i-1}$ then, for any value of $(x_0, x_1, \ldots, x_{i-1})$, $\omega_i(x_0, x_1, \ldots, x_i) = R_i$ for some outcome $x_i \in [\underline{x}_i, \bar{x}^i]$ . Furthermore the payment scheme has a flat segment at the reference and, therefore, $\omega_i$ is not be strictly increasing.*

3. *If $\omega_{T-1}(x_0, x_1, \ldots, x_{T-1}) \leq R_{T-1}$ then, for any value of $(x_0, x_1, \ldots, x_{T-1})$, $\omega_T(x_0, x_1, \ldots, x_T) = R_T$ for some outcome $x_T \in [\underline{x}_T, \bar{x}^T]$. The payment scheme has a flat segment at the reference but it cannot be flat in $x_T$.*

*Proof.* The payment scheme must be continuous and non-decreasing, in fact. The multipliers $(\mu_i)_i$ are strictly positive (same al canonical model). It can be seen that the

right side of (1.26) has slope with respect to $x_i$ of at least $\frac{d}{dx_i}\left((1-\delta\ell_{i+1})\mu_i\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}\right)$ and therefore, $\omega_i$ must be non-decreasing. A flat segment will arise whenever the payment scheme reaches the first period reference level. (1.20) and (1.21) imply that, as $x_i$ increases and reaches the reference if the scheme were to continue increasing it would enter the gain area and the right side of (1.20) and (1.21) would jump downwards, therefore contradicting that it increased after reaching the reference income. Something analogous happens when the reference level is reached from above (as $x_i$ decreases), if the optimal scheme were to go below the reference, the optimality characterization would require it to jump upwards. This contradicts that it decreased after reaching the reference from above.

Now, let's see whether the reference will be reached. The following equality must be fulfilled,

$$\frac{1}{U'(\omega_i(x_0,x_1,\ldots,x_i))} = (1+k_i(x_0,x_1,\ldots,x_i)\ell_i)\left(\lambda_i+\mu_i\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}\right)+$$
$$-\delta\ell_{i+1}\int k_{i+1}(x_0,x_1,\ldots,x_{i+1})(\lambda_{i+1}+\mu_{i+1}\frac{f^{i+1}_{a_{i+1}}(x_{i+1}|a_{i+1})}{f^{i+1}(x_{i+1}|a_{i+1})})f^i(x_{i+1}|a_{i+1})dx_{i+1}. \quad (1.14)$$

The $i-1$ period's payments fulfills the following equation,

$$\lambda_i(1+k_{i-1}(x_0,x_1,\ldots,x_{i-1})\ell_{i-1}) = \frac{1}{U'(\omega_{i-1}(x_0,x_1,\ldots,x_{i-1}))}+$$
$$\delta\ell_i\int k_i(x_0,\ldots,x_i)\left(\lambda_i+\mu_i\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}\right)f^i(x_i|a_i)dx_i.$$

Suppose $\omega_i(x_0,x_1,\ldots,x_i) < \omega_{i-1}(x_0,x_1,\ldots,x_{i-1}) \quad \forall x_i$ (except in one point). Then $k_i(x_0,x_1,\ldots,x_i)=1 \quad \forall x_i$ and (1.15) becomes $\lambda_i = \frac{1}{U'(\omega_{i-1})(1+k_{i-1}\ell_{i-1}-\delta\ell_i)}$. Therefore, since $\delta\ell_{i+1}\int k_{i+1}(x_0,x_1,\ldots,x_{i+1})(\lambda_{i+1}+\mu_{i+1}\frac{f^i_{a_{i+1}}(x_{i+1}|a_{i+1})}{f^i(x_{i+1}|a_{i+1})})f^i(x_{i+1}|a_{i+1})dx_{i+1} \le \delta\ell_{i+1}\lambda_{i+1}$ we obtain

$$\frac{1}{U'(\omega_i)} \ge (1+\ell_i-\delta\ell_{i+1})\left(\frac{1}{U'(\omega_{i-1})(1+k_{i-1}\ell_{i-1}-\delta\ell_i)}+\mu_i\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}\right)$$

Therefore, if $(\ell_i-\ell_{i+1})\delta \ge k_{i-1}\ell_{i-1}-\ell_i$ and $\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)} > 0$ we conclude $\frac{1}{U'(\omega_i)} \ge \frac{1}{U'(\omega_{i-1})}$ which contradicts $\omega_i < \omega_{i-1}$. If $\ell_{i-1}=\ell_i$, $(\ell_i-\ell_{i+1})\delta \ge \ell_{i-1}-\ell_i$ or $k_{i-1}=0$ then $(\ell_i-\ell_{i+1})\delta \ge k_{i-1}\ell_{i-1}-\ell_i$ will be fulfilled.

Suppose $\omega_i(x_0,x_1,\ldots,x_i) > \omega_{i-1}(x_0,x_1,\ldots,x_{i-1}) \quad \forall x_i$ (except in one point). Then

27

$k_i(x_0, x_1, \ldots, x_i) = 0 \quad \forall x_i$ and (1.15) becomes $\lambda_i = \frac{1}{U'(\omega_{i-1})(1+k_{i-1}\ell_{i-1})}$. Therefore, we obtain

$$\frac{1}{U'(\omega_i)} \leq \left( \frac{1}{U'(\omega_{i-1})(1 + k_{i-1}\ell_{i-1})} + \mu_i \frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)} \right)$$

Therefore, if $\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)} < 0$ we conclude $\frac{1}{U'(\omega_i)} \leq \frac{1}{U'(\omega_{i-1})}$ for which contradicts $\omega_i > \omega_{i-1}$. We conclude that the reference must be reached on an interval for all the cases stated above.

$\square$

The previous property states that the payment scheme will also be non-decreasing and that the reference point must be paid for some outcome. Furthermore, the reference must be paid for an interval in the outcomes support. If we assume a continuous distribution function, the only payment the agent will receive with non-zero probability will be the reference. The size of the flat segment will depend on the parameters of the model. Figure 1.3.1.1 is an schematic representation of the monotonicity of possible payment schemes. According to Property 4 (a) is possible in periods $\{0, \ldots, T-1\}$, (b) is possible only in period 0, and (c), (d) and (e) are possible in all periods.

There are two properties that distinguish the shape of optimal payment scheme from the classical case. First, the optimal payment scheme can have flat segments. This is explained by the multiplicative term $(1 + k_i(x_0, x_1, \ldots, x_i)\ell)$ in 1.20 and 1.21. At the reference level of the agent, $k_i(x_0, x_1, \ldots, x_i)$ is allowed to take any value in $[0, 1]$. Therefore, the right hand side of 1.20 and 1.21 can remain constant in an interval, as $x_i$ increases, $k_i(x_0, x_1, \ldots, x_i)$ decreases and $\omega_i(x_0, x_1, \ldots, x_i)$ remains at the reference level $R$.

The second difference with the classical relates to the fact that the Principal takes into account that the each period's payment affects the reference level of the following period. The last term of the right side of 1.20 represents this effect. It is strictly positive if the $i + 1$th period payment scheme has an interval of results for which the payment is lower or equal than the reference of the agent. From Property 4 we know that this will be the case in every period. This term will tend to lower the payment scheme and
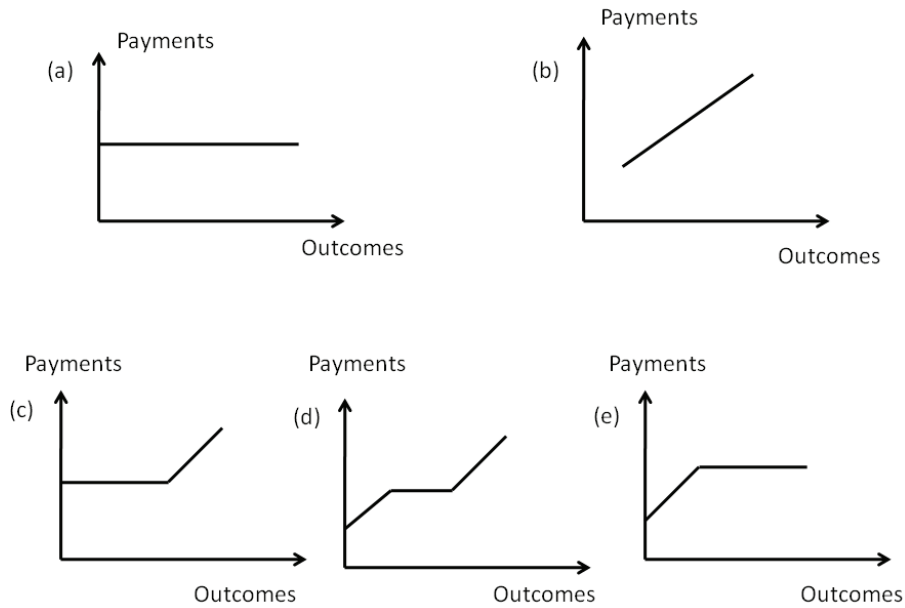
28

Figure 1.2: Schematic representation of monotonicity of contracts

reduce its growth rate as $x_i$ grows. It represents the benefit of a lower reference level when the $i + 1$th wage schedule pays below the reference for some outcomes, which is a higher utility function of the agent.

This result is analogous to De Meza and Webb (2007). The existence of flat segments on the reference level of the payment schedule is optimal and is a consequence of a kinked utility function. For an outcome that pays the reference level to the agent, it might be costly to give a payment over or under the reference for close results. If the payment goes slightly below the reference the agent will have a relatively high loss in utility, because of the higher marginal utility below the reference, thus straining (PC). Similarly, if the payment goes slightly over the reference the increased incentive will be low because the agent will experience a relatively high fall in the marginal utility for payments over the reference level.

**Inter-temporal Properties**

Just as in the classical case, consumption smoothing requires that a higher payment

29

in one period implies, ceteris paribus, a higher payment in all subsequent periods. As stated earlier, in our model in every period two different outcomes may pay the agent's reference level. This implies, ceteris paribus, that the wage schedules in later periods are greater or equal (as functions) and they may overlap for outcomes that pay the reference level. In the classical model the wage schedules are strictly increasing as functions. This is discussion is summarized in the following property.

**Property 5** (Dependence across periods). *Let $x'_i < x''_i$ two possible outcomes in period $i$.*

1. *If $\omega_i(x_0, x_1, \ldots x'_i) = \omega_i(x_0, x_1, \ldots x''_i) = R_i$ then*

$$\omega_j(x_0, x_1, \ldots x'_i, x_{i+1}, \ldots, x_j) \leq \omega_j(x_0, x_1, \ldots x''_i, x_{i+1}, \ldots, x_j) \quad \forall j > i \quad \forall x_j \in [\underline{x}_j, \overline{x}_j]$$

   .

2. *If $\omega_i(x_0, x_1, \ldots x'_i) < \omega_i(x_0, x_1, \ldots x''_i)$ then*

$$\omega_j(x_0, x_1, \ldots x'_i, x_{i+1}, \ldots, x_j) < \omega_j(x_0, x_1, \ldots x''_i, x_{i+1}, \ldots, x_j) \quad \forall j > i \quad \forall x_j \in [\underline{x}_j, \overline{x}_j]$$

   .

*Proof.* It follows directly from 1.20 and 1.21 since $\mu_i > 0 \quad \forall i$

$\square$

Property 5 implies that pictures (a) and (b) in Figure 1.3.1.1 are possible in our model. Note also that only (a) is possible in the canonical model.

In Rogerson (1985) a relationship between the wage schedules of two consecutive periods. The inverse of the marginal utility of income must equal the conditional expected value of the inverse of the marginal utility. In mathematical terms, the following must be fulfilled,

$$\frac{1}{u'(Y_{i-1}(x_0, x_1, \ldots, x_{i-1}))} = \int_{\underline{s}_1}^{\overline{s}_1} \frac{1}{u'(Y_i(x_0, x_1, \ldots, x_i))} f^i(x_i|a_i) dx_i \qquad (1.15)$$

where $u(\cdot)$ is a differentiable utility function, $Y_{i-1}(x_0, x_1, \ldots, x_{i-1})$ and $Y_i(x_0, x_1, \ldots, x_i)$ are the optimal payment schemes of period's $i-1$ and period $i$ respectively.
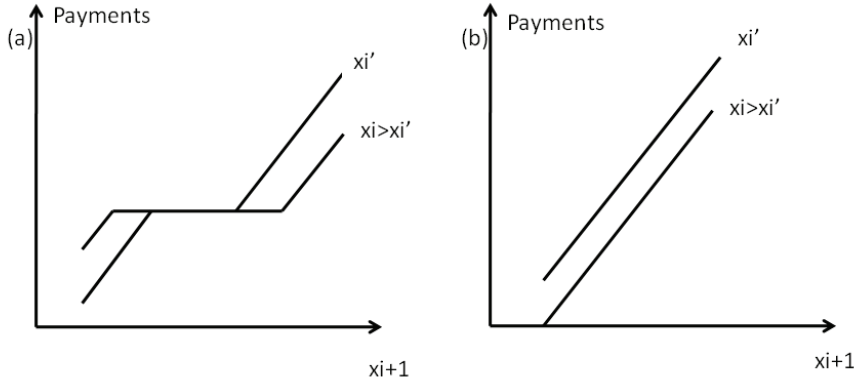
30

Figure 1.3: Schematic representation of monotonicity of contracts

This condition is no longer valid in our model. However, an extended condition can be derived as is stated in the following property,

**Property 6** (Two consecutive periods relationship)**.** *The following relationship between two consecutive periods must be fulfilled,*

$$\frac{1}{U'(\omega_{i-1}(x_{i-1}))(1 + k_{i-1}(x_{i-1})\ell_{i-1})} = \int \frac{1}{U'(\omega_i(x_i))(1 + k_i(x_i)\ell_i)} f^i(x_i|a_i)dx_i + c(x_{i-1})$$

*where*

$$
c(x_{i-1}) = -\frac{\ell_i \delta}{1 + k_{i-1}(x_{i-1})\ell_{i-1}} \int k_i(x_i) \left( \lambda_i + \mu_i \frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)} \right) f^i(x_i|a_i)dx_i + \tag{1.16}
$$
$$
\ell_{i+1} \delta \int \int \frac{k_{i+1}(x_{i+1})}{1 + k_i(x_i)\ell_i} \left( \lambda_{i+1} + \mu_{i+1} \frac{f^{i+1}_{a_{i+1}}(x_{i+1}|a_{i+1})}{f^{i+1}(x_{i+1}|a_{i+1})} \right) f^{i+1}(x_{i+1}|a_{i+1})f^i(x_i|a_i)dx_{i+1}dx_i
$$

9

*Proof.* Follows directly from 1.20 and 1.21.

$\square$

---

[9]Some outcome variables have been omitted for simplicity.

Property 6 implies that the inverse of the marginal utility of income might be greater or smaller than than the conditional expectation of the marginal utility. If the payment is always strictly over the reference in period $i+1$ and period $i+2$ the classical equality maintains.[10]

In the canonical model, consumption smoothing and incentive provision requires the principal to give relatively higher payments in each period with respect to the next one. This facilitates the provision of incentives in later periods. Therefore, after each period, if allowed to save or borrow, the agent has incentives to save. In our model, this will not be the case, the agent may have incentives to save, borrow or consume his allocation. Furthermore, he faces a somewhat higher incentive to consume the allocation that he was given. To see this, suppose the loss averse agent of our model faces the possibility of allocating resources between two consecutive periods, $i$ and $i+1$. If the payment scheme in period $i+1$ pays the reference in a set of results with measure greater than 0 he faces the following conundrum. If he ponders borrowing he takes into account the increased utility in the current period and the decrease in utility in the following period, when he returns what he was borrowed. The increase in utility in the current period will depend whether his allocation is over or under the reference. It will be relatively higher strictly under the reference. The decrease in utility in period $i+1$ will be high on and under the reference, which by assumption is a set with positive probability. Furthermore, the reference in period $i+1$ is increased thus reducing the utility of that period. This effects may make borrowing unattractive to the agent. For instance, if the agent were paid the reference in one period the increase in utility will be low with respect to the decrease in utility of paying back when receiving a payment equal to the reference in the following period. Similarly, when deciding whether or not to save he realizes that under the reference the loss in utility would be high with respect to the gain in utility on or over the reference in the following period, which can make saving unattractive. This implies that the agent will face situations in which he would face a loss in utility by saving or by borrowing, and decides to consume the allocation that he was given.

---

[10]There is some abuse of language here since we refer as marginal utility of $\tilde{U}_i$ to the term $U'(\omega_i(x_i))(1 + k_i(x_i)\ell_i)$ since both quantities are equal for all incomes except the reference. At the reference level the marginal utility is not computable and could take any value between $[U'(\omega_i(x_i)), U'(\omega_i(x_i))(1 + \ell_i)]$.

The previous discussion describes a result that can be interpreted as a a "status quo bias", as was first described in Samuelson and Zeckhauser (1988). Imagine the agent is in a situation in which he would lose if the decides either to save or to borrow at interest rate $1/\delta - 1$[11]. If facing the possibility of lending a part of his income at interest rate $r_l$ and borrowing at interest rate $r_b$, indifference between lending and borrowing, i.e. the marginal utility of borrowing and savings is zero, requires $r_l > \frac{1}{\delta} - 1 > r_b$. This means that the smallest price at which he is willing to lend part of his income is strictly greater than the greatest price at which he is willing to borrow. This is not the case in the canonical model since, indifference between lending and borrowing, for any set of payments schemes for two consecutive periods, can be attained for a single interest rate[12]. This gap between willingness to accept and willingness to buy has been described in other economical contexts (Bateman et al. (1997)). If the market has interest rates such that $r_b \geq r_l$, the agent will find himself inclined to consume his allocation. A formalization of the previous discussion is presented in the following property.

**Property 7** (Status quo bias).    ■ *If the payment in period $i$ is $y$ and is at the reference and the payment in period $i + 1$ is constant such that $y_{i+1}(x_{i+1}) = y \quad \forall x_{i+1} \in [\underline{x}_{i+1}, \bar{x}_{i+1}]$, then the agent will neither want to save nor borrow at rate $\frac{1}{\delta} - 1$ and if $r_l$ is the rate that makes marginal utility of saving equal to zero and $r_b$ the rate that makes the marginal utility of borrowing is equal to zero, we must have $r_l > \frac{1}{\delta} - 1 > r_b$.*

     ■ *Let $y_i$ be the payment in period $i$ and $y_{i+1}(x_{i+1})$ the payment scheme in period $i+1$. If $y_{i+1}$ pays the reference with positive probability or $y_i$ is at the reference, then the rate that makes marginal utility of saving equal to zero $r_l$ is strictly greater than the rate that makes the marginal utility of borrowing, $r_b$, equal to zero.*

*Proof.* The marginal utility of saving in period $i$ at rate $r_l$ and consuming the savings

---

[11]Under the assumption that both principal and agent have discount factor $\delta$, $\frac{1}{\delta} - 1$ is the market's interest rate.

[12]As long as we let the interest rate be negative.

in period $i+1$ is given by,

$$-(1+1_{\{\omega_i \le R_i\}}\ell_i)U'(\omega_i) \quad + \quad \delta(1+r_l)\int (1+\ell_{i+1}1_{\{\omega_{i+1}<\omega_i\}})U'(\omega_{i+1}(x_{i+1}))f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1} \quad (1.17)$$

$$+\ell_{i+1}\delta U'(\omega_i)\int_{\omega_{i+1}<\omega_i} f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1}$$

and the marginal utility of borrowing in period $i$ at rate $r$ and paying back in period $i+1$ is given by,

$$(1+1_{\{\omega_i < R_i\}}\ell_i)U'(\omega_i) \quad - \quad \delta(1+r_b)\int (1+\ell_{i+1}1_{\{\omega_{i+1}\le\omega_i\}})U'(\omega_{i+1}(x_{i+1}))f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1} \quad (1.18)$$

$$-\ell_{i+1}\delta U'(\omega_i)\int_{\omega_{i+1}\le\omega_i} f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1}$$

By (1.17) the marginal utility of saving at rate $\frac{1}{\delta}-1$ is $-(1+\ell)U'(y)+U'(y) < 0$. By (1.18) the marginal utility of borrowing at rate $\delta$ is $U'(y)-(1+\ell_{i+1})U'(y)-\delta\ell_{i+1}U'(y) < 0$. Therefore the rate at which the agent would be willing to borrow is smaller than $\frac{1}{\delta}-1$ and the rate at which he would be willing to save must be greater than $\frac{1}{\delta}-1$.

The second point is justified subtracting (1.17) and (1.18) with $r_l = r_b = r$ and noting that what is obtained is strictly negative.

□

A more technical explanation for the previous property is that if the agent is allowed to save or borrow after one period he will decide to save if the marginal utility of savings is strictly positive and to borrow if the marginal utility of borrowing is strictly positive. The loss averse preferences of our model imply that in each period the marginal utility of savings does not coincide with the minus the marginal utility of borrowing and, therefore, they may be both negative simultaneously.

Under the optimal payment scheme the agent may have incentives to save or to borrow in our model. There are several effects at stake. First the ones described above, that relate to the the fact that marginal utility varies depending on whether the payments are under or over the reference and the nature of the set of outcomes for which the reference is paid. Also to the fact that while saving or borrowing the agent influences his own reference for the following period. The optimality of the payment scheme will also affect the willingness to save or borrow. The Principal wants to decrease each period's payment in order to increase the agent's utility in the following period, thus favoring

34

borrowing in the current period. The agent will want to increase savings in order to increase utility in the second period, thus giving incentives to save. Finally, in order to facilitate the provision of incentives in the future, just as in the classical case, the optimal payment may be higher than the agent would like thus giving incentives to save.

The following property analyzes whether the agent would like to save or borrow at interest rate $\frac{1}{\delta} - 1$ if allowed to do so under some possible shapes of the payment schemes in our model.

**Property 8** (Inter-temporal allocation of resources)**.** *If the interest rate is $\frac{1}{\delta} - 1$ in period $i$ then the agent may have incentives to save for period $i + 1$, to borrow and pay back in period $i + 1$ or to consume exactly its income depending on the parameters of the problem. Moreover,*

- *If period's $i + 1$ payment scheme is over the reference for all results, then the agent does not have incentives to save in period $i$.*

- *If period's $i + 1$ payment scheme is below the reference for all outcomes in period $i + 1$ then*

  - *if period's $i$ payment is strictly above the reference then the agent has incentives to save in period $i$.*

  - *If period's $i$ payment is at the reference, the agent will not have incentives to borrow in period $i$.*

  - *If period's $i$ payment is strictly below the reference the agent may have incentives to save, to borrow or to consume his allocation.*

*Proof.* Suppose that $\omega_{i+1}(x_{i+1}) \geq \omega_i$ for all $x_{i+1}$. By (1.17) the marginal utility of saving would be

$$-(1 + 1_{\{\omega_i \leq R_i\}}\ell_i)U'(\omega_i) + \int U'(\omega_{i+1}(x_{i+1}))f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1} \tag{1.19}$$

By assumption we will have that $\int U'(\omega_{i+1}(x_{i+1}))f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1} < U'(\omega_i)$ and therefore the agent will not have incentives to save. By 1.18 the marginal utility of borrowing is

$$(1 + 1_{\{\omega_i < R_i\}}\ell_i)U'(\omega_i) - \int (1 + \ell_{i+1}1_{\{\omega_{i+1}=\omega_i\}})U'(\omega_{i+1}(x_{i+1}))f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1}$$

$$-\ell_{i+1}\delta U'(\omega_i)\int_{\omega_{i+1}=\omega_i} f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1}$$

it may be positive or negative depending on the parameters of the problem

Now, suppose that $\omega_{i+1}(x_{i+1}) \leq \omega_i$ for all $x_{i+1}$. The marginal utility of saving is be

$$-(1+1_{\{\omega_i \leq R_i\}}\ell_i)U'(\omega_i) + \int (1+\ell_{i+1}1_{\{\omega_{i+1}<\omega_i\}})U'(\omega_{i+1}(x_{i+1}))f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1}$$
$$+\ell_{i+1}\delta U'(\omega_i)\int_{\omega_{i+1}<\omega_i} f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1}$$

We must have $\int U'(\omega_{i+1}(x_{i+1}))f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1} > U'(\omega_i)$, therefore, if $\omega_i > R_i$ or $\ell_i$ is sufficiently small then the marginal utility would be positive and therefore the agent will have incentives to save. The marginal utility of borrowing is,

$$(1+1_{\{\omega_i < R_i\}}\ell_i)U'(\omega_i) - \int (1+\ell_{i+1}1_{\{\omega_{i+1}\leq\omega_i\}})U'(\omega_{i+1}(x_{i+1}))f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1}$$
$$-\ell_{i+1}\delta U'(\omega_i)\int_{\omega_{i+1}\leq\omega_i} f_{i+1}(x_{i+1}|a_{i+1})dx_{i+1}$$

$\square$

The previous property states, among other things, that if the payment scheme is flat for lower outcomes and then increasing for greater outcomes, then the agent will not have incentives to save. Later in this section we show that observing a payment scheme that does not venture under the reference is possible as long as the cost to the principal of the participation constraint $\lambda$ is sufficiently high.

## 1.3.2.  Renegotiation-proofness and spot-implementability

The optimal contract scheme will be renegotiation-proof just as in the classical case. This is because we are assuming that the agent is able to predict how her utility function updates in each period. If this weren't the case the renegotiation-proofness could be broken. [13]

However, the optimal sequence of spot contracts will exhibit memory, unlike the classical case, because the assumption we make on the dynamical update of the reference level. The utility function of the agent changes from period to period, and given that we make this assumption the reservation utility from period to period may change too. However the optimal commitment contract may not be implemented by spot contracts.

---

[13]For a reference note that the proof the renegotiation-proofness property in Chiappori et al. (1994) does not rely on the differentiability of the utility function.

Suppose that the reservation utility of the agent in each period $i$ is given by $\bar{U}(R_i)$. In summary we will have,

**Property 9** (Optimal spot contracting). *The optimal sequence of spot contracts exhibits memory and it will not implement the full-commitment optimum in general.*

*Proof.* By backwards induction, the optimal spot contract in period $T$ must give the agent the reservation utility $\bar{U}(\omega_{T-1})$ and will depend on $\omega_{T-1}$ since it represents the reference in period $T$. Thus, the optimal sequence of spot contracts has memory. The optimal spot contract in period $T-1$ solves

$$\max_{\omega_{T-1}(\cdot)} \int \left( (x_{T-1} - \omega_{T-1}(x_{T-1})) f^{T-1}(x_{T-1}|a_{T-1}) + \delta V(\omega_{T-1}(x_{T-1})) f(x_T|a_T) \right) dx_{T-1}$$

$$\int \left( \tilde{U}_{T-1}(\omega_{T-1}(x_{T-1})) + \delta \bar{U}(\omega_{T-1}(x_{T-1})) \right) f^{T-1}(x_{T-1}|a_{T-1}) \geq \bar{U}(c_{T-2})$$

$$a_{T-1} \in \operatorname{argmax}_a \int \left( \tilde{U}_{T-1}(\omega_{T-1}(x_{T-1})) + \delta \bar{U}(\omega_{T-1}(x_{T-1})) \right) f^{T-1}(x_{T-1}|a_{T-1})$$

where $V(\omega_{T-1}(x_{T-1}))$ represents the profits of the principal under the optimal spot contract in period $T$. Therefore, unless $\bar{U}(\omega_{T-1}(x_{T-1}))$ coincides with the expectation of the last period contract under the full-commitment optimum the optimal sequence of spot contracts does not implement the full-commitment optimum.

$\square$

### 1.3.3. The shape of the optimal contract and the multiplier $\lambda$

In a three period context we analyze how the optimal payment schemes changes if we change the multiplier $\lambda$ which represents the cost for the principal of the participation constraint (PC). We have the following property,

**Property 10** (Cost of (PC) and shape of optimal contract). *There is a value $\bar{\lambda}$ such that the if $\lambda \geq \bar{\lambda}$ the optimal payment scheme is over the reference in all three periods (letting the other multipliers be fixed).*

*Proof.* See appendix

$\square$

This property can be extended to any number of periods. It suggests that the cost of the participation constraint to the principal is closely related to whether incentives are to be created through rewards or punishments. This result is highly intuitive. Providing incentives through punishments in the loss area creates a great loss in utility for the

agent whose PC is already very costly. The principal is better off generating incentives by rewards only.

## 1.4.  Monitorable access to credit

If there is monitorable access to credit, savings can be contracted upon. The program that the principal faces is the following,

$$\max_{(\omega_i(\cdot))_i, (a_i)_i, (s_i)_i, (S_i)_i} \sum_{t=0}^{T} \delta^t \mathbb{E}\left(x_t - \omega_t(x_0, x_1, \ldots, x_t) - S_t(x_0, x_1, \ldots, x_t) | a_0, a_1, \ldots, a_t\right)$$

subject to

$$\sum_{t=0}^{T} \delta^t \left(\mathbb{E}\left(\tilde{U}_t(\omega_t(x_0, x_1, \ldots, x_t) - s_t(x_0, x_1, \ldots, x_t), c_{t-1}) | a_0, a_1, \ldots, a_t\right) - \psi_i(a_i)\right) \geq U^* \tag{PC}$$

$$(a_0, a_1(x_1), \ldots a_T(x_0, x_1, \ldots, x_T)) \in \text{argmax}_{\vec{a}} \sum_{t=0}^{T} \delta^t \left(\mathbb{E}\left(\tilde{U}_t(\omega_t(x_0, x_1, \ldots, x_t) - s_t(x_0, x_1, \ldots, x_t), c_{t-1}) | a_0, a_1, \ldots, a_t\right) - \psi(a_i)\right)$$

$$\tag{IC}$$

where $s_t$ are the agent's accumulated savings in period $t$. That is the net savings of the agent in period $t$ once the endowment derived from previous savings is taken into account. $S_t$ are the accumulated savings of the principal. It is easy to see that the previous program is equivalent to one in which the optimization variables are the consumptions of the agent in each period, given by $c_t(x_0, x_1, \ldots, x_t) = \omega_t(x_0, x_1, \ldots, x_t) - s_t(x_0, x_1, \ldots, x_t)$, and the total accumulated savings, given by $s_t + S_t$, and constraint $s_T = -\frac{s_{T-1}}{\delta}$. Since we are assuming that the principal is risk neutral the optimality conditions will be similar to (1.20) and (1.21) with $c_i$ replacing the payments $\omega_i$. Therefore, the optimality conditions for consumptions with monitorable access to credit will be the following,

$$\frac{1}{U'(c_i(x_0, x_1, \ldots, x_i))} = (1 + k_i(x_0, x_1, \ldots, x_i)\ell_i)\left(\lambda_i + \mu_i \frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}\right) +$$

$$- \delta\ell_{i+1} \int k_{i+1}(x_0, x_1, \ldots, x_{i+1})(\lambda_{i+1} + \mu_{i+1} \frac{f^i_{a_{i+1}}(x_{i+1}|a_{i+1})}{f^i(x_{i+1}|a_{i+1})})f^i(x_{i+1}|a_{i+1})dx_{i+1}. \qquad \forall i < T \quad (1.20)$$

and

$$\frac{1}{U'(c_T(x_0, x_1, \ldots, x_T))} = (1 + k_T(x_0, x_1, \ldots, x_T)\ell_T)\left(\lambda_T + \mu_T \frac{f^T_{a_T}(x_T|e_i)}{f^i(x_i|e_i)}\right) \tag{1.21}$$

where $k_i(x_0, x_1, \ldots, x_i)$ is given by,

$$k_i(x_0, x_1, \ldots, x_i) \in \begin{cases} \{1\} & \text{if } \omega_i(x_0, x_1, \ldots, x_i) < R_i \\ [0,1] & \text{if } \omega_i(x_0, x_1, \ldots, x_i) = R_i \\ \{0\} & \text{otherwise} \end{cases}$$

This result is similar to what is obtained in the classical case. There is a strong relationship between the monitorable access to credit case and the full-commitment with no credit access case. Furthermore, just as in the classical case, monitoring borrowing and savings introduces memory to the principal-agent relationship and therefore the optimal long-term contract will be spot-contractible.

**Property 11** (Spot contractibility under monitorable credit). *Suppose the reservation utility $U_i^*(s_{i-1}, R_i)$ in period $i$ depends on the savings that the agent has in period $i-1$ and the reference level $R_i$ and that it is continuous, increasing in $s_{i-1}$ then, under monitorable savings the long-term optimal contract will be spot contractible.*

*Proof.* By backward induction, if in the last period the contract requires consumptions $c_T(x_0, x_1, \ldots, x_T)$ and $s_{T-1}$ is such that the reservation utility in period $T$ equals the agent's expected utility under the optimal contract. That is, if $U^*(s_{T-1}, c_{T-1})$ is the reservation utility of the agent with accumulated savings $s_{T-1}$ and reference $c_{T-1} = \omega_{T-1} - s_{T-1}$ and the following equality is fulfilled,

$$\mathbb{E}(\tilde{U}_T(c(x_T)) | a_T) - \psi(a_T) = U^*(s_{T-1}, c_{T-1})$$

Then, since the long-term contract is ex-post efficient, the last period spot contract will implement the last period contract of the optimal long-term contract. Note that under the assumptions of continuity of $U^*$ and unlimited transfers there will be a value for $s_{T-1}$ that will satisfy the equality given the consumption plan $c(x_T)$ and effort level $a_T$.

Now in period $T-1$ the principal knowing $(x_T, s_T, a_T)$, will be accepted in the following period, offers contract $(x_{T-1}, s_{T-1}, a_{T-1})$ which is spot contractible if $s_{T-2}$ is such that $U^*(s_{T-2}, c_{T-2}) = \mathbb{E}(\tilde{U}(c(x_{T-1})) + \delta\tilde{U}(c(x_T)) | a_{T-1}, a_T) - \psi(a_{T-1}) - \delta\psi(a_T)$. This contract will be accepted by the agent since he knows what spot contract he will be offered in period $T$ and is optimal for the principal by ex-post efficiency of optimal long-term contract. Inductively one concludes that the long-term contract can be implemented by spot-contracting.

$\square$

## 1.5. Two period example

In this section we compute numerically the optimal payment schemes in a two period setting in order to illustrate the possible shapes the optimal contracts can take. The distribution function of outcomes $x_i \in [0,1]$ in period $i \in \{1,2\}$, for effort level $a_j \in \{a_L, a_H\}$ is a triangular given by

$$f^i(x_i|a_j) = \begin{cases} \frac{2x_i}{a_j} & x_i \leq a_j \\ \frac{2(1-x_i)}{1-a_j} & x_i > a_j \end{cases} \tag{1.22}$$

and $U(Y) = \sqrt{Y}$ and therefore

$$\widetilde{U}_i(Y_i, R_i) = \sqrt{Y_i} - \theta(Y_i, R_i)\ell_i(\sqrt{R_i} - \sqrt{Y_i})$$

Note that in this case $\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}$ may not be strictly increasing with respect to outcomes $x_i$. In what follows we assume $a_H = 1$ and $a_L = 0.1$ in which case $\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}$ is constant in $[0, 0.1]$.

The optimality conditions are,

$$\frac{1}{U'(\omega_0(x_0))} = 2\sqrt{\omega_0(x_0)} = (1 + k_0(x_0)\ell_0)\left(\lambda_0 + \mu_0 \frac{f^0_{a_0}(x_0|a_0)}{f^0(x_0|a_0)}\right) +$$
$$- \delta\ell_1 \int_{\omega_1 \leq \omega_0} k_1(x_0, x_1)(\lambda_1 + \mu_1 \frac{f^1_{a_1}(x_1|a_1)}{f^1(x_1|a_1)})f^1(x_1|a_1)dx_1. \tag{1.23}$$
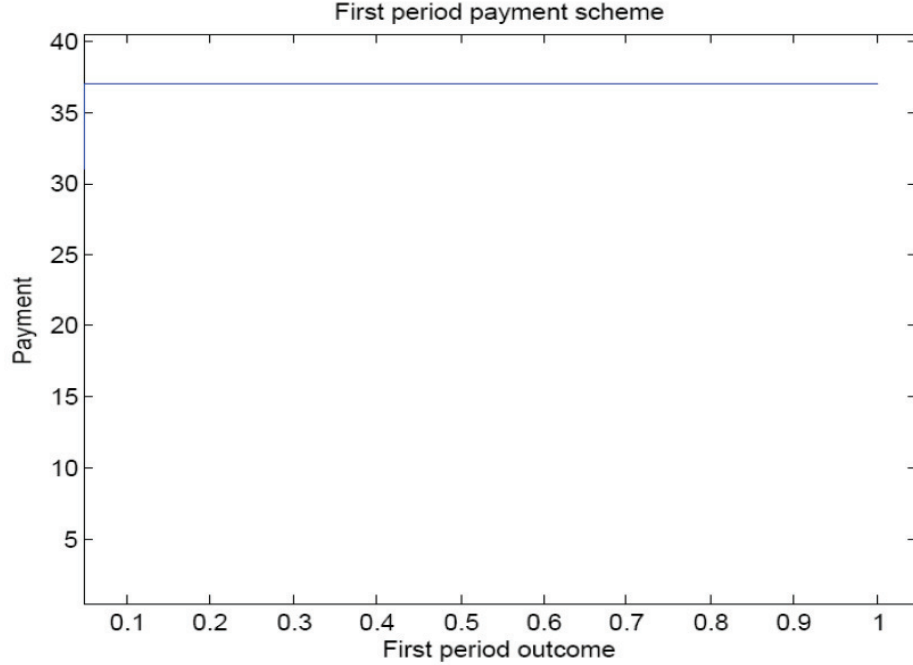
$$\frac{1}{U'(\omega_1(x_0,x_1))} = 2\sqrt{\omega_1(x_0,x_1)} = (1 + k_1(x_0,x_1)\ell_1)\left(\lambda_1 + \mu_1 \frac{f^1_{a_1}(x_1|a_1)}{f^1(x_1|a_1)}\right) \tag{1.24}$$

### 1.5.1. Case 1: First period payment independent of outcomes

The first example illustrates a case in which the first period payment does not depend on the outcomes that take place in first period (see Figure 1.5.1). That is the first period payment stays constant at the reference level. The second period payment scheme is contingent on outcomes obtained on the first and second period and can be seen in Figure 1.5.1.[14] The values of parameters used in this simulation are given in

---

[14]Note that the flat segments for small values of first and second period outcomes is due to $\frac{f^i_{a_i}(x_i|a_i)}{f^i(x_i|a_i)}$ constant in $[0, 0.1]$

Figure 1.4: Case 1. First Period Payment Scheme



the following table,

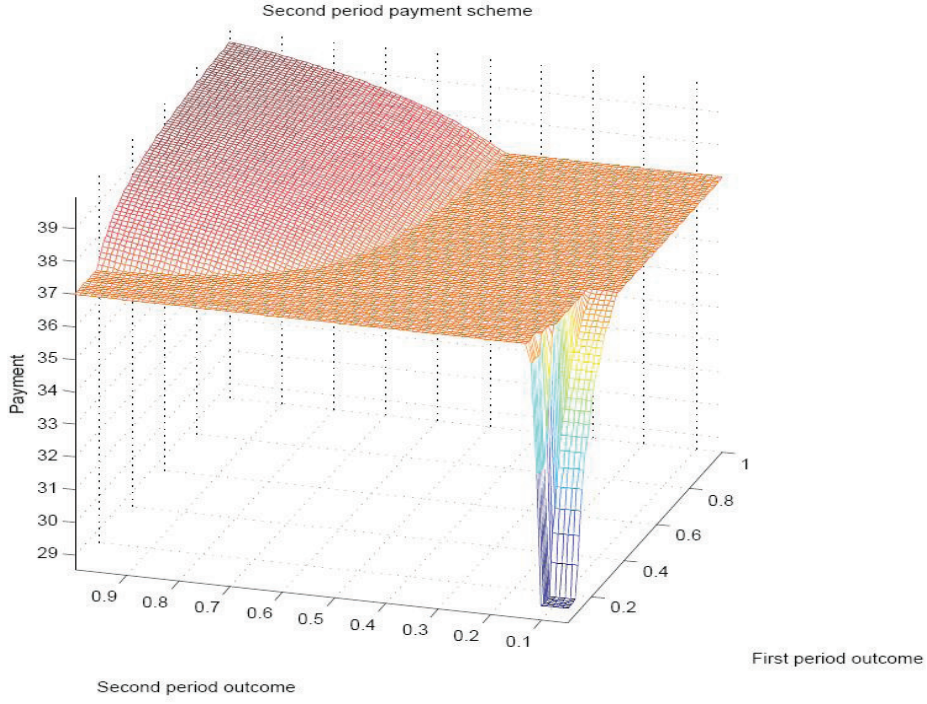| $\lambda$ | $\mu_0$ | $\mu_1$ | $\ell_0$ | $\ell_1$ | $a_H$ | $a_L$ | $1/U'(R_0)$ |
|---|---|---|---|---|---|---|---|
| 46.1 | 0.5 | 2 | 1 | 1 | 1 | 0.1 | 37 |

Note that the second period payments fall below the reference for the low outcomes in the first period. For outcomes in the first period that are greater than a threshold, the agent will face a payment scheme in the following period that is greater or equal than the payment received in the first period. Consequently, according to Property 8, he will not have incentives to save for outcomes above a threshold.

## 1.5.2.   Case 2: First period payment greater or equal than $R_0$

The next example illustrates a case in which the first period payment reaches the reference for low values of the first period outcome (see Figure 1.5.2). The second period payment scheme is shown in Figure 1.5.2. The values of parameters used in this simulation are given in the following table,

| $\lambda$ | $\mu_0$ | $\mu_1$ | $\ell_0$ | $\ell_1$ | $a_H$ | $a_L$ | $1/U'(R_0)$ |
|---|---|---|---|---|---|---|---|
| 30.1 | 1 | 1 | 1 | 1 | 1 | 0.1 | 15 |

41

Figure 1.5: Case 1. Second Period Payment Scheme



## 1.5.3. Case 3: Second period payment over reference for all outcomes

This example illustrates a case in which payments are over the reference for all outcomes in the first and second period (see Figure 1.5.3). The parameters used are the same as in Case 2, except that $\lambda$ is greater, thus illustrating Property 10. The second period payment scheme is contingent on outcomes obtained on the first and second period and can be seen in Figure 1.5.3. According to Property 8, the agent does not have incentives to save for any outcome in the first period. The payment scheme shown in Figures 1.5.3 and 1.5.3 is efficient, renegotiation-proof and implements the high level of effort in both periods if the agent is restricted to borrow. The values of parameters used in this simulation are given in the following table,

| $\lambda$ | $\mu_0$ | $\mu_1$ | $\ell_0$ | $\ell_1$ | $a_H$ | $a_L$ | $1/U'(R_0)$ |
|---|---|---|---|---|---|---|---|
| 40.1 | 1 | 1 | 1 | 1 | 1 | 0.1 | 15 |

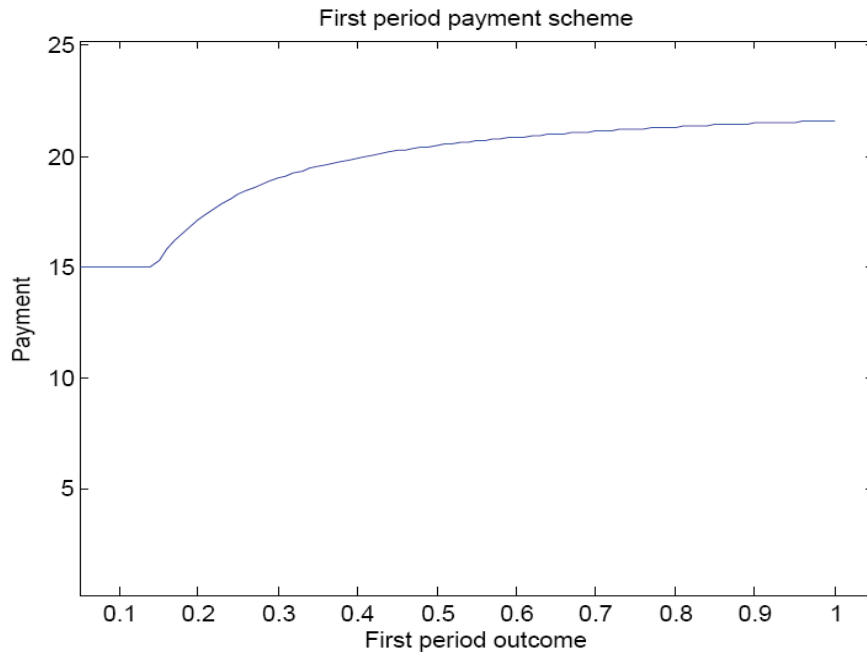Figure 1.6: Case 2. First Period Payment Scheme



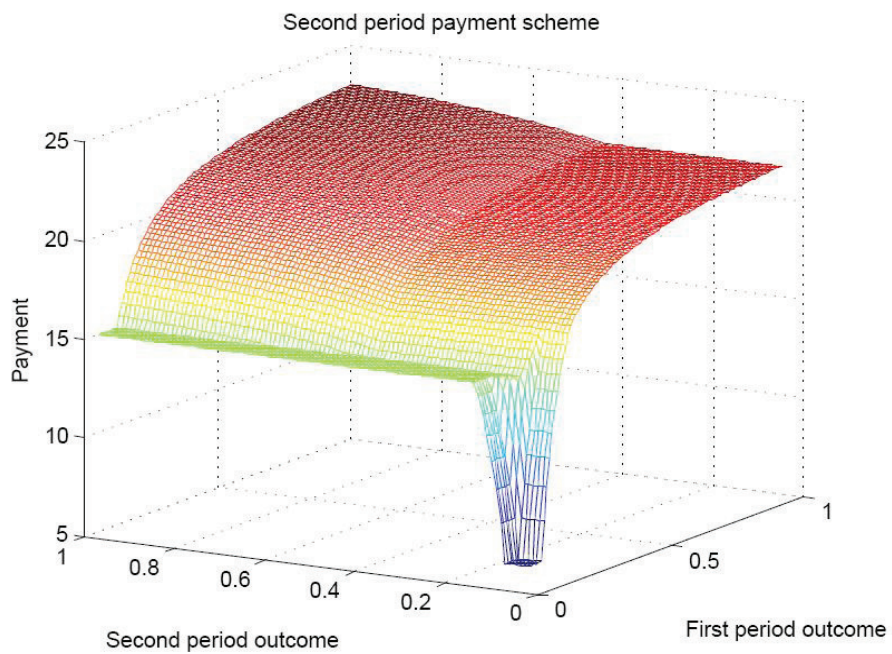Figure 1.7: Case 2. Second Period Payment Scheme

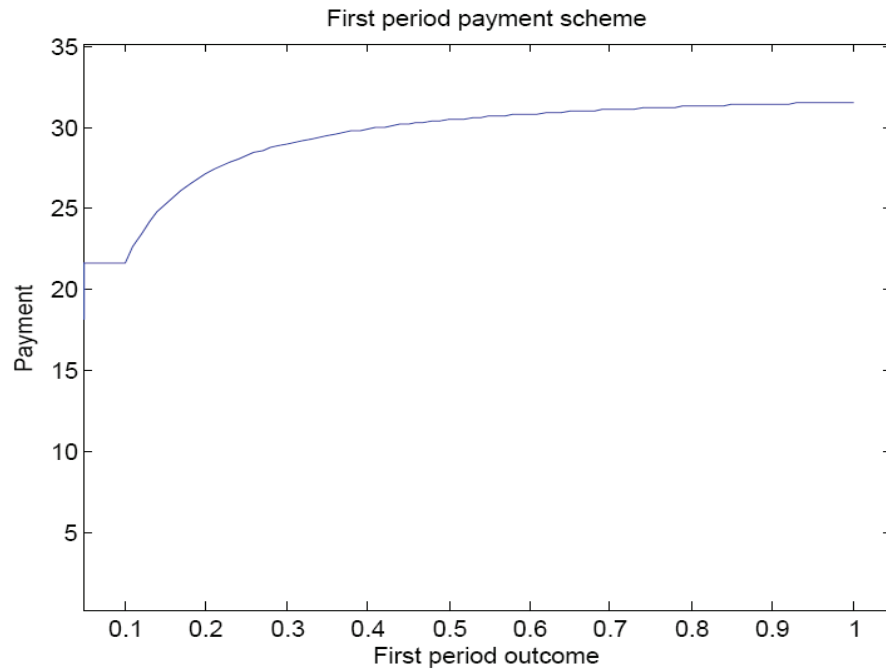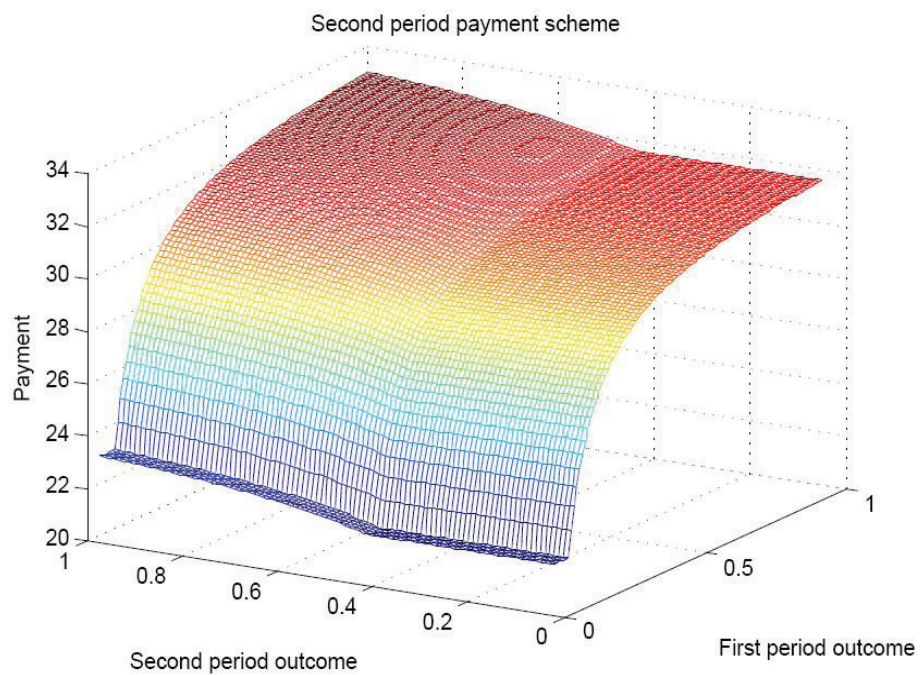Figure 1.8: Case 1. First Period Payment Scheme



Figure 1.9: Case 1. Second Period Payment Scheme

# 1.6.  Conclusions and Final Remarks

We introduce loss aversion to the canonical dynamic moral hazard model introduced by Rogerson (1985). We make the assumption that the reference level corresponds to the consumption that took place in the previous period. We find differences and similitudes in terms of predictions with the canonical model. In particular, we find that payment schemes may have flat segments in outcomes in the first period and must have flat segments from the second period on. The flat segment may extend for the entire support of the outcomes distribution. We find that in comparing with the canonical model, payments schemes will exhibit a smaller dependence in outcomes within and across periods.

We derive an extended relationship between the optimal payment schemes of two consecutive periods. We find that the result of the canonical model that requires that after each period the agent is inclined to save some of his earnings may not be valid depending on the parameters of the model. This, together with ex-post efficiency of the full-commitment optimum, implies that efficiency and high incentive provision may be attained under no savings constraints, or no credit constraints altogether. The latter being justified by the presence of a "status quo bias" in the inter-temporal preferences defined by our model, that we identify as a gap between the lowest interest rate at which the agent is willing to save and the highest interest rate at which he is willing to borrow.

By introducing loss aversion into a dynamic principal-agent model with endogenous update of the reference level, we provide an explanation for a number of deviations of observed contracts from the theoretical predictions of the canonical model. In particular, it allows us to explain the existence of contracts that are not strictly increasing in outcomes and the weak history dependence of contracts found in the empirical literature.

Finally, the derivation of the optimality conditions is not exempt of mathematical difficulty given that the objective function is non-differentiable. subdifferential analysis is required to find the solution. A "chain rule" not previously described in subdifferential calculus that applies to the problem at hand is derived. As a consequence, this paper also represents a methodological guide to tackle this or other similar problems that may involve loss aversion or not. It is important to note that in many contexts, the

results and insights of models might substantially change once the assumption of non-differentiability is included, since the usual first order conditions are no longer valid.

## 1.7. Appendix

### 1.7.1. Proof of property 10

Using the fact that for any period $i$ at the reference $R_i$ (i.e. when $k_i \in [0,1]$) we must have

$$k_i(x_0, x_1, \ldots, x_i)\left(\lambda_i + \mu_i \frac{f_{a_i}^i(x_i|a_i)}{f^i(x_i|a_i)}\right) = \frac{1}{U'(R_i)} - \left(\lambda_i + \mu_i \frac{f_{a_i}^i(x_i|a_i)}{f^i(x_i|a_i)}\right)$$

This implies that the three period optimality conditions are the following,

$$\frac{1}{U'(\omega_0(x_0))}\left(1 + \delta\mathbb{P}(\omega_1 = \omega_0(x_0)|a_1) + \delta^2\mathbb{P}(\omega_2 = \omega_0(x_0)|a_2)\mathbb{P}(\omega_1 = \omega_0(x_0)|a_1)\right) =$$

$$(1 + k_0(x_0)\ell_0 + \delta\mathbb{P}(\omega_1 = \omega_0(x_0)|a_1) + \delta^2\mathbb{P}(\omega_2 = \omega_0(x_0)|a_2)\mathbb{P}(\omega_1 = \omega_0(x_0)|a_1) +$$

$$- \delta\ell_1\mathbb{P}(\omega_0(x_0) > \omega_1|a_1) - \delta^2\ell_2\mathbb{P}(\omega_0(x_0) > \omega_1|a_1)\mathbb{P}(\omega_0(x_0) > \omega_2|a_2))\left(\lambda + \mu_0\frac{f_{a_0}^0(x_0|a_0)}{f^0(x_0|a_0)}\right) +$$

$$+ -\delta^2\ell_2\int_{\omega_1=\omega_0,\omega_2<\omega_0}\left(\mu_1\frac{f_{a_1}^1(x_1|a_1)}{f^1(x_1|a_1)} + \mu_2\frac{f_{a_2}^2(x_2|a_2)}{f^2(x_2|a_2)}\right)f_{a_1}^1(x_1|a_1)f_{a_2}^2(x_2|a_2)dx_1dx_2 +$$

$$+ \delta^2\int_{\omega_1=\omega_0,\omega_2=\omega_0}\left(\mu_1\frac{f_{a_1}^1(x_1|a_1)}{f^1(x_1|a_1)} + \mu_2\frac{f_{a_2}^2(x_2|a_2)}{f^2(x_2|a_2)}\right)f_{a_1}^1(x_1|a_1)f_{a_2}^2(x_2|a_2)dx_1dx_2$$

$$- \delta\mu_1\ell_1\int_{\omega_0<\omega_1}f_{a_1}(x_1|a_1)dx_1 + \mu_1\delta\int_{\omega_0=\omega_1}f_{a_1}(x_1|a_1)dx_1 \quad (1.25)$$

$$\frac{1}{U'(\omega_1(x_0,x_1))}\left(1 + \delta\mathbb{P}(\omega_2 = \omega_1(x_0,x_1)|a_2)\right) = (1 + k_1(x_1)\ell_1 + \delta\mathbb{P}(\omega_1(x_0,x_1) = \omega_2|a_2) +$$

$$- \delta\ell_1\mathbb{P}(\omega_1(x_0,x_1) > \omega_2))\left(\lambda_1 + \mu_1\frac{f_{a_1}^1(x_1|a_1)}{f^1(x_1|a_1)}\right)$$

$$+$$

$$- \delta\mu_2\ell_2\int_{\omega_2<\omega_1}f_{a_2}(x_2|a_2)dx_2 + \mu_1\delta\int_{\omega_2=\omega_1}f_{a_1}(x_1|a_1)dx_1 \quad (1.26)$$

$$\frac{1}{U'(\omega_2(x_0,x_1,,x_2))} = (1 + k_2(x_0,x_1,x_2)\ell_2)\left(\lambda_2 + \mu_2\frac{f_{a_2}^2(x_2|a_2)}{f^2(x_2|a_2)}\right) \quad (1.27)$$

46

From (1.25) it can be seen that for whatever value of the first period's reference a $\lambda$ high enough will insure that right hand side is greater than $1/U'(R_0)$ for $k_0 = 0$. Now, replacing in (1.26) and (1.27) the following,

$$\left(\lambda + \mu_0 \frac{f^0_{a_0}(x_0|a_0)}{f^0(x_0|a_0)}\right) = \quad \frac{1}{U'(\omega_0(x_0))} +$$

$$-\delta^2 \int_{\omega_1=\omega_0,\omega_2=\omega_0} \left(\mu_1 \frac{f^1_{a_1}(x_1|a_1)}{f^1(x_1|a_1)} + \mu_2 \frac{f^2_{a_2}(x_2|a_2)}{f^2(x_2|a_2)}\right) f^1_{a_1}(x_1|a_1) f^2_{a_2}(x_2|a_2) dx_1 dx_2 +$$

$$-\mu_1 \delta \int_{\omega_0=\omega_1} f_{a_1}(x_1|a_1) dx_1$$

we conclude that if $1/U'(\omega_0(x_0)$ is big enough the right hand sides of both equations will be greater than $1/U'(\omega_0(x_0))$ which implies that the payment schemes will be over the reference in periods 1 and 2.

# Bibliography

I. Bateman, A. Munro, B. Rhodes, C. Starmer, R. Sugden, 1997. A Test of the Theory of Reference-Dependent Preferences. Quarterly Journal of Economics, May 1997, Vol. 112, No. 2, pages 479-505.

Bertsekas D.P., Nedic A., Ozdaglar A.E. Convex analysis and optimization. Athena Scientific, 2003. ISBN 1886529450.

P. Bolton, M. Dewatripont, 2005. Contract Theory. The MIT Press.

D. Bowman, D. Minehart, M. Rabin, 1999. Loss Aversion in a Consumption-Savings Model. Journal of Economic Behavior & Organization, February 1999, Vol. 38, pages 155-178.

C. D. Carroll, M. S. Kimball, 2001. Liquidity Constraints and Precautionary Saving. NBER Working Paper No. W8496, August 2001.

P. Chiappori, I. Macho, P. Rey and B. Salanié. 1994. Repeated Moral Hazard: The role of memory, commitment, and the access to credit markets. European Economic Review, Elsevier, October 1994, Vol. 38, pages 1527-1553.

P.A. Chiappori, B. Salanié, 2000. Testing Contract Theory: a Survey of Some Recent Work. Invited lecture World Congress of the Econometric Society Seattle, August 2000.

De Meza, D. Webb, 2007. Incentive Design under Loss Aversion. Journal of the European Economic Association, MIT Press, March 2007, Vol. 5, pages 66-92.

F. Gul, 1991. A Theory of Disappointment Aversion. Econometrica, May 1991, Vol. 59, No. 3, pages 667-686.

SJ. Grossman, OD. Hart, 1983. An analysis of the principal-agent problem. Econometrica, January 1983, Vol. 51, pages 7-45.

GW Harrison, JA List, 2004. Field experiments. Journal of Economic Literature, 2004, Vol. 42, pages 1009-1045.

B. Hölmostrom, 1979. Moral hazard and observability. The Bell Journal of Economics, Spring 1979, Vol. 10, No. 1, pages 74-91.

B. Hölmstrom; P. Milgrom, 1987. Aggregation And Linearity In The Provision Of Intertemporal Incentives. Econometrica, March 1987, Vol. 55, pages 303-328.

D Kahneman, A Tversky, 1979. Prospect theory: An analysis of decision under risk. Econometrica, March 1979, Vol. 47, No. 2, pages 263-292.

B. Köszegi, M. Rabin, 2006. A Model of Reference-Dependent Preferences. Quarterly Journal of Economics, November 2006, Vol. 121, No. 4, pages 1133-1165.

A. Munro, R. Sugden, 2003. On the theory of reference dependent preferences. Journal of Economic Behavior and Organization, April 2003, Vol. 50, pages 407-428.

Canice Prendergast, 2002. The Tenuous Trade-Off between Risk and Incentives. The Journal of Political Economy, October 2002, Vol. 110, No. 5, pages 1071-1102.

William P. Rogerson, 1985. Repeated Moral Hazard. Econometrica, January 1985, Vol. 53, No. 1, pages 69-76.

W. Samuelson, R. Zeckhauser, 1988. Status Quo Bias in Decision Making. Journal of Risk and Uncertainty, March 1988, Vol. 1, No. 1, pages 7-59.

S. Shavell, 1979. Risk Sharing and Incentives in the Principal and Agent. The Bell Journal of Economics, Spring 1979, Vol. 10, No. 1, pages 55-73.

Stole, Lars, 2001. Lectures on the Theory of Contracts and Organizations.

R. Tyrrel Rockafellar, 1974. Conjugate Duality and Optimization. CBMS-NSF Regional Conference Series In Applied Mathematics.

# Chapter 2

# Contract Theory applied to On-demand IT Services Contracting

Nicolas Figueroa      Alejandro Jofre      Sofía Moroni

Akhil Sahai      Yuan Chen      Subu Iyer

## 2.1. Introduction

An Service Level Agreement (SLA/Contract) is an agreement between a provider and a consumer which is comprised of Service Level Objectives that guarantee quality of service (such as availability, performance and reliability), a promise of payment and penalties to impose in case the objectives are not met. The study of such contracts has become increasingly important with the increasing use of IT outsourcing procedures, which had reached $56 billion in 2000 and was expected to reach $100 by 2005 (Dermikan et al. (2005)). While the original practice of IT outsourcing contracts involved complicated measures to safeguard the client's interest against the many potential mishaps, a more modern approach has focused on a system of penalties and rewards based on observed quality of service, serving as a monetary compensation that insures the client in case the service is suboptimal (Dermikan et al. (2005)).

In this work we focus on the problem of offering optimal (revenue maximizing) contracts from the Service Providers' (SP) point of view. In particular, we are interested in contracts offered by IT providers, that offer service guarantees in terms of performance, availability, security and reliability constraints. These contracts specify the pricing for

the service guarantees and the penalties that are due in case of violations. We model SLA/Contracts using the concepts of Moral Hazard and Adverse Selection.

The Moral Hazard comes from the fact that the provider, through some costly effort (investment, use of scarce resources such as number of CPU's, number of engineer hours, etc.), can increase the quality of the service, but that there is also an additional stochastic component to it. The effort level cannot be monitored by the client, and the actual performance of the system (that the client can observe) is just a noisy signal of effort. In an IT context, better infrastructure on average provides better performance, but some unforeseen incidents (unforeseen demand increase, breakdown of a system, etc.) may still lead to poor quality. Since effort is not observable, the only way to induce a high level of effort is through a compensation system that is "steep" i.e. with higher payments when observed quality is better, or equivalently with penalties if the providers does not meet his end of the deal. Nonetheless, this affects the provider, since she may sometimes be punished for low quality even if the effort put in the process was high. Given her risk aversion, she will demand higher expected payments when the payment system is steeper. The basic trade-off is then set: "steeper" compensation systems will induce higher effort, but they will shift more risk (in terms of earnings) to the provider, who is risk averse and will charge more for the service. We introduce then the "credibility constraints": a contract must promise a level of effort that is optimal given the penalties imposed in case of non-compliance with the quality level promised. Any other effort level would not be credible and the client would not accept such a contract.

At the same time, the service provider is faced with an adverse selection problem: clients differ in their valuation of the service, in their risk aversion and in other characteristics. Moreover, their particular characteristics are private information, and the service provider only knows the distribution of possible clients. In order to deal with this issue, the service provider must offer different contracts (one for each type of client he can potentially face) and design them in such a way that clients choose the contract that was designed for them. Such constraints (called the selfselection constraints) decrease the revenue an SP may obtain from clients, since they extract an informational rent due to asymmetric information. We construct a general model incorporating both the credibility and self-selection constraints, and allowing for risk averse clients and service provider. Since we are interested in the practical application of such a model

to the case of a service provider in the IT sector, we allow for a general shape of the stochastic relation between effort and quality, and proceed to numerically solve for the optimal pricing policy. This optimal policy includes different contracts (tailored to be selected by the different types of clients), each one specifying a fixed payment and a bonus based on the quality delivered.

The literature on SLA/Contracts has addressed many issues: the type of quality measurements that can be used to define SLA agreements and the type of verification that must be used to satisfy both parties (M.J. Buco et al. (2004)), the necessity of imposing penalties as a way to insure the client against bad performance (Dermikan et al. (2005)) and the optimal software design for a SP serving multiple clients and allocating scarce resources (K. Appleby et al. (2001)). However, the literature has been silent on two critical issues. First, the use of penalties to induce credible levels of effort from the SP since effort is non observable. Second, and even more important: the optimal design of such penalties, as a function of the client and SP's risk aversion, and the stochastic relationship between effort and quality. There has also been work on optimal pricing, in the case of a monopolist who faces clients with private information (Chone et al. (2001)), but this issue has not addressed the penalties issue, since it is assumed that clients are "small" and are provided with a generic quality level for which they cannot complain. That approach is correct for many markets (for example cellphones, where a client is never compensated when the network does not perform as expected), but is clearly unrealistic in the case of IT outsourcing, where the SP and the client are of similar size and have similar bargaining power.

The paper is organized as follows. In section two we present a single-client model and an extension to a multi-client scenario. Section tree discusses some practical questions for finding an SLA/Contract using our approach. Section four and five analyzes SLA determination in a n-tier IT service scenario, including numerical results.

## 2.2. Mathematical Model

In this section we present a basic model to price SLA/Contracts for IT services, in which there is one client and a service provider, and an extension in which there are different types of clients. In both cases, we let the contract depend on the realized

quality of service. This contract, in turn, must give incentives to the service provider to put in the effort level (infrastructure, servers, labor, etc) that is implicitly agreed upon.

## 2.2.1. Basic Model

The provider delivers a quality of service $q_{\in}Q \subseteq \mathbb{R}^{n_q}$ to the client. This quality level can be throughput, availability, response time, etc, and depends stochastically on the effort level $e \in E \subseteq \mathbb{R}^{n_e}$ that the provider assigns to the client. In this setting, clients are not able to observe the level of effort (for example, servers which could internally allocated to another task )that the SP has assigned to each one of them. We assume that the distribution function of quality, given a level of effort $e$ is $f_q(q|e)$ with cumulative function $F_q(q|e)$, and its support will be a cube in $\mathbb{R}^{n_q}$, $\Pi_{p=1}^{n_q}[\underline{q_p}, \bar{q}_p]$, independent of $e$. We assume that it is possible to transform the observable and verifiable variable $q$ into a quality of service in monetary terms $q_m$. This transformation is represented by a function $g : Q \to Q_M$ which is non-decreasing and concave in each component of $q$.

Since the level of effort is not observable by the client, a contract can only specify a payment contingent on quality, that is observable and verifiable. We assume that the payment will be a function of monetary quality of service, which is observable since it is a deterministic function of quality. We will denote this payment rule as $p(q_m)$. As we will see there is no loss of generality in considering contracts that depend on $q_m$ and not on $q$. Notice that $e$ being non-observable introduces a "moral hazard" problem, the provider must carry out a hidden action, which is beneficial to the client, and non-contractible.

We denote by $V$ the utility function of a client. For a realized level of monetary quality $q_m$, his utility will be $V(q_m - p(q_m))$. On the other hand, we assume that the provider has a utility function that depends on effort and money, so for a given level of effort $e$ and quality $q_m$, her utility is $U(p(q_m), e)$. As usual, we assume $V' > 0$ and $V'' \leq 0$ and that $\frac{\partial U}{\partial p} > 0$, $\frac{\partial U}{\partial e} < 0$, $\frac{\partial^2 U}{\partial^2 p} \leq 0$ and $\frac{\partial^2 U}{\partial^2 e} \geq 0$ (see figure 2.1). This assumptions imply that both the service provider and the client are risk-averse.

Under these assumptions, and based on the tools of contract theory, we are able to state the problem of finding a utility maximizing contract for the service provider as a constrained optimization problem.

There are two constraints. The *participation constraint* (2.2 from now on "PC"), states that the utility level of the client has to be above a certain exogenous level $\bar{V}$, reflecting the opportunity cost of the resources involved. The *credibility constraint* (2.3, from now on "CC") states that the promised effort level has to be optimal (for the service provider) given the contract, since otherwise the client would not trust her.

Therefore, the provider chooses $e$ and $p(q_m)$ to solve

$$\max_{p(),e} \int U(p(q_m), e) f_{q_m}(q_m|e) dq_m \tag{2.1}$$

subject to

$$\int V(q_m - p(q_m)) f_{q_m}(q_m|e) \geq \bar{V} \tag{2.2}$$

$$e \in \text{argmax}_{e'} \int U(p(q_m), e') f_{q_m}(q_m|e') dq_m \tag{2.3}$$

For future use, we also write the problem in terms of $q$:

$$\max_{p(),e} \int U(p(g(q)), e) f_q(q|e) dq \tag{2.4}$$

subject to

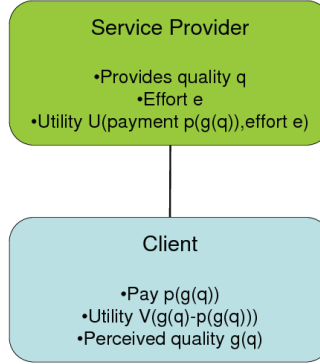$$\int V(g(q) - p(g(q))) f_q(q|e) dq \geq \bar{V} \tag{2.5}$$

$$e \in \text{argmax}_{e'} \int U(p(g(q)), e') f_q(q|e') dq \tag{2.6}$$

## 2.2.2. The General Model

We can enrich the model to take into account that the SP could have different clients, that can be classified in "types". The agents that belong to a particular type will differ from the other types because they can have different valuations of the service, and utility functions. The Service Provider will offer a menu of contracts, to satisfy the different necessities of the clients and therefore extract more payments.

The main constraint is that Service Provider does not know which client is which, since the particular valuation for the service is private information of the client, or that he has
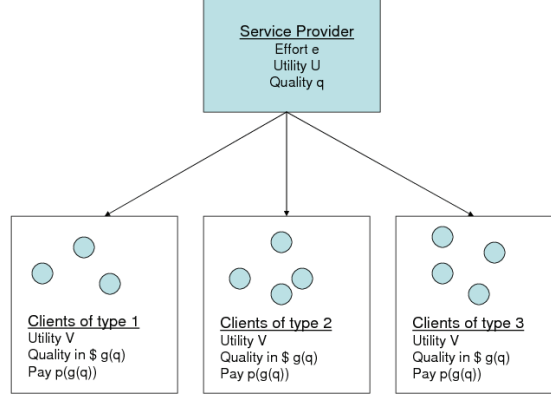
to offer all contracts to all clients for legal reasons. Given this constraint, the contracts have to satisfy a *self-selection* property: each client selects the contract that was designed for his particular type because is the one that gives him the greatest utility level.

For the model, let's suppose that there are $N$ clients, which can be classified in $k$ types. We denote by $\theta_j$ the "type" of agents in the class $j$, and we assume that there is an amount $\mu_j$ of those agents. The Service Provider will offer a menu of contracts, each of them designed for one type of client. As before the contracts will offer a payment given the realized monetary quality $q_m^i$ that client $i$ receives. We will denote the contract intended for agents of type $\theta_j$ as $p_{\theta_j}(q_m)$. The Service Provider will devote independent efforts for each client, which we will denote $\{e_{i,\theta_j}\}_{i\in\{1,\ldots\mu_j\},j\in\{1\ldots k\}}$, that is, the effort that the i'th client of type $\theta_j$ will be assigned, will be $e_{i,\theta_j}$. The utility function of the Service Provider will take the form $U(\{p_{\theta_j}(\cdot)\}_j, \sum_i e_i)$, that is, it will depend on payments and the sum of efforts dedicated to each client, with $\frac{\partial U}{\partial p_i} > 0$, $\frac{\partial U}{\partial e_i} < 0$, $\frac{\partial^2 U}{\partial^2 p_i} \leq 0$ and $\frac{\partial^2 U}{\partial^2 e_i} \leq 0$ (see figure 2.2)

Clients of the same type, say $\theta_j$ will have the same utility function $V(\cdot|\theta_j)$, the same $g_{\theta_j}$ function and reservation utility $\bar{V}(\theta_j)$, that will relate true quality to monetary quality of service, and distribution function of quality $q_{\theta_j} \in \Pi_{p=1}^{n_{q_{\theta_j}}} [\underline{q}_{p,\theta_j}, \bar{q}_{p,\theta_j}] \subset \mathbb{R}^{n_{\theta_j}}$ [1] contingent on effort, $f_{\theta_j}(q_{\theta_j}|e)$. The distribution function of monetary quality will depend on the effort level that each client receives, and will be independent of the effort given to other clients. The optimization problem that the Service Provider will face is,

---

[1]Notation analogous to Basic Model

Figure 2.2: Multi-Client Model

$$\max_{\{p_{\theta_j}(\cdot)\}_j,\{e_{i,\theta_j}\}_{j,i\in\{1,\dots,\mu_j\}}} \int U(\{p_{\theta_j}(g_{\theta_j}(q^{i,\theta_j}))\}_{i,j}, \sum_{i,j} e_{i,\theta_j}) \prod_{j,i} f_{\theta_j}(q^{i,\theta_j}|e_{i,\theta_j}) d\overrightarrow{q} \qquad (2.7)$$

subject to

$$\int V(g_{\theta_i}(q) - p_{\theta_i}(g_{\theta_i}(q))|\theta_i) f_{\theta_i}(q|e_{r,\theta_i}) dq \geq \int V(g_{\theta_i}(q) - p_{\theta_j}(g_{\theta_j}(q))|\theta_i) f_{\theta_i}(q|e_{t,\theta_j}) dq \qquad \forall j \neq i, \forall r,t$$
$$(2.8)$$

$$\int V(g_{\theta_i}(q) - p_{\theta_i}(g_{\theta_i}(q))|\theta_i) f(q|e_i) \geq \bar{V}(\theta_i) \qquad (2.9)$$

$$\{e_{i,\theta_j}\}_{j,i\in\{1,\dots,\mu_j\}} \in \mathrm{argmax}_{\{e'_{i,\theta_j}\}} \int U(\{p_{\theta_j}(g_{\theta_j}(q^{i,\theta_j}))\}_{i,j}, \sum_{i,j} e'_{i,\theta_j}) \prod_{j,i} f_{\theta_j}(q^{i,\theta_j}|e'_{i,\theta_j})$$
$$(2.10)$$

As before the objective is to maximize the expected utility of the service provider. (2.8) is the *Self Selection Constraint*, it states that in the optimal contract each client type will prefer their own contract to the ones intended for other client types. (2.9) and (2.10) are the *Participation Constraint* and *Credibility Constraint*, respectively.

The optimization problem presented above is complicated. If we introduce a number of assumptions on the utility functions and the distributions of quality given the efforts the framework will be much simplified (details in appendix).

For simplicity and in order to illustrate the economic intuitions of the results, in this work we focus in optimizing over functions $p(\cdot)$ that are linear. That is, we are looking for contracts that are linear in monetary quality of service. This contracts are simple and could easily be implemented in real applicatons. With the assumptions made (2.7)-(2.10) and it will become,

$$\max_{\{p_{\theta_j}(\cdot)\}_j, \{e_{\theta_j}\}_j} \int U(\{p_{\theta_j}(g_{\theta_j}(q^{\theta_j}))\}_j, \sum_j \mu_j e_{\theta_j}) \prod_j f_{\theta_j}(q^{\theta_j}|e_{\theta_j}) \tag{2.11}$$

subject to

$$\int V(g_{\theta_i}(q) - p_{\theta_i}(g_{\theta_i}(q))|\theta_i) f_{\theta_i}(q|e_{\theta_i})dq \geq \int V(g_{\theta_i}(q) - p_{\theta_j}(g_{\theta_j}(q))|\theta_i) f_{\theta_i}(q|e_{\theta_j})dq \quad \forall j \neq i \tag{2.12}$$

$$\int V(g_{\theta_i}(q) - p_i(g_{\theta_i}(q))|\theta_i) f(q|e_i) \geq \bar{V}(\theta_i) \tag{2.13}$$

$$\{e_{\theta_j}\}_j \in \mathrm{argmax}_{\{e'_{\theta_j}\}} \int U(\{p_{\theta_j}(g_{\theta_j}(q^{\theta_j}))\}_j, \sum_j \mu_j e'_{\theta_j}) \prod_j f_{\theta_j}(q^{\theta_j}|e'_{\theta_j}) \tag{2.14}$$

That is, all clients of the same type will have the same amount of effort assigned and the expected utility of the provider will be simplified to an integral over $\sum_{j=1}^{k} n_{q_{\theta_j}}$ variables. In the framework above we assumed that the utility of the SP and the client could be computed as an expected utility. Now we want to allow for the possibility that the utility of the agents doesn't have that form, that is, it is not a Bernoulli Utility Function. Our framework will remain the same, except that now the utilities of the agents will not be written as an expected utility, and they will be a function of the random variable of profits, that will, in turn, depend on the level of effort. If $X$ is a random variable that describes the behavior of uncertain profits, with pdf $f_X(x)$ and cdf $F_X(x)$ an example of a function of the sort is what we refer to the expectation-variance utility function given by the formula $U(X) = I\!\!E(X) - \tau \int (x - I\!\!E(X))^2 f_X(x)dx$, with $\tau$ a non-negative constant. The first term is the expected profits and the second term is the variance of profits. Since the variance is a measure of how volatile profits can be, the second term implies that agents utility decreases with risk. Another example is what we will the CVaR utility, $U(X) = I\!\!E(X) - \tau(-I\!\!E(X_\alpha))$ where $X_\alpha$ is the r.v. known as the lower $\alpha$-tail of $X$, with distribution function $F_{X_\alpha} = \frac{min\{\alpha, F_X\}}{\alpha}$, $\alpha \in [0, 1]$, and $\tau$ a non-negative constant. $-I\!\!E(X_\alpha)$ is known as the CVaR and it is a measure of risk. For instance, if $F_X(\cdot)$ is continuous, the CVaR will be the mean of the lowest $\alpha\%$

of the profits.

## 2.3.  Optimal Contract

Solving (2.4)-(2.6) is, in general, difficult, because of the last constraint. However, under certain conditions, if the maximization problem in (2.6) has an interior solution which is its unique stationary point[2], the last constraint can be simplified to an equality constraint which corresponds to the First Order Condition of the Optimization Problem in (2.3).[3] This will happen, in particular, when the expected utility of the SP is concave in effort at the optimal contract.

Replacing the (2.6), in the basic model, or (2.10), in the multi-client case, by a first order condition is known as the first order approach (FOA). The conditions under which it is valid will depend on $U$ and $f_q(q|e)$. In this section we study conditions that allow us, in the context of SLAs, to apply such a method.

### 2.3.1.  First Order Approach

A sufficient condition to be able to use the FOA is that the objective function be concave in $e$. If $q$ is uni-dimensional we have the following results. From Jewitt (1988),

**Lemma 1.** *If $F_q(q|e)$ satisfies (2.15)-(2.16), then, for every function $\tilde{u}(\cdot) : \mathbb{R} \to \mathbb{R}$ in $C^1$ concave and non-decreasing, $\int \tilde{u}(q) f_q(q|e) dq$, will be concave in $e$ and therefore, if $U(p(q_m), e) = u(p(q_m)) - \phi(e)$, with $u$ concave and $\phi(\cdot)$ convex in $e$, then the optimization problem (2.1)-(2.3), solved for $p(\cdot)$ of linear form with positive slope, will satisfy FOA.*

$$\int_{-\infty}^{y} F_q(q, e) dq \text{ is nonincreasing convex in } e \text{ for each value of } y \qquad (2.15)$$

$$\int_{-\infty}^{\infty} q f_q(q, e) dq \text{ is nondecreasing and concave in } e \qquad (2.16)$$

**Corollary 1.** *If $F_q(q|e)$ is convex in $e$ for each $q$, then, if $U(p(q_m), e) = u(p(q_m)) - \phi(e)$ with $u$ concave and $\phi(\cdot)$ convex in $e$, $\int u(p(g(q))) f_q(g(q)|e) dq$, will be concave in $e$, and therefore, the optimization problem (2.1)-(2.3), solved for $p(\cdot)$ of linear form with positive slope, will satisfy FOA.*

---

[2] An stationary point is one in which the gradient of the SP's utility with respect to $e$ becomes 0

[3] Note that this applies also for the case in which the set of efforts is not $\mathbb{R}^{n_e}$: the first order condition will correspond to the gradient of a Lagrangian.

This results can be extended to the case in which we have multi-dimensional $q = (q_1, q_2, \ldots, q_{n_q})$, if there are separability conditions of the utility function $U(p(g(q)), e)$ in the components of the vector $q$. This also applies to the multi-client case. Conditions (2.15)-(2.16) can guarantee the validity of the first order approach in some contexts. Therefore it is useful to study the distributions $f(x|e)$ that will satisfy them. Next we present some examples of contexts in which they will be satisfied.

**Example 1.** *Let $F(q, \mu, \sigma)$ be the cdf of the normal distribution with mean $\mu(e)$ and variance $\sigma^2(e)$. If $R(\frac{y-\mu}{\sigma}) \cdot \sigma$ is convex and non-increasing in effort for every $R$ such that $R' \geq 0$ and $R'' \geq 0$, then $F(x, \mu, \sigma)$ will satisfy condition (2.15). The same applies for the log-normal distribution. A proof of this result can be found in the appendix. A similar result can be derived for a truncated normal. In fact if $q$ has pdf,*

$$f_q(q|e) = \begin{cases} \dfrac{\frac{e^{-\frac{(t-\mu)^2}{\sigma}}}{\sigma\sqrt{2\pi}}}{d(e)} & \text{if } t \in [-c\sigma + \mu, c\sigma + \mu] \\ 0 & \sim \end{cases}$$

*Where $d(e) = \int_{-r(e)-\mu}^{r(e)-\mu} \frac{e^{\frac{(x+\mu)^2}{2\sigma^2}}}{\sigma\sqrt{2\pi}} = \int_{-\frac{r(e)}{\sigma}}^{\frac{r(e)}{\sigma}} \frac{e^{\frac{x^2}{2}}}{\sqrt{2\pi}}$, $\mu(e)$ and $\sigma(e)$ depend on $e$ and $c$ is a constant.*

*If $R(\frac{y-\mu}{\sigma}) \cdot \sigma$ is convex and non-increasing in effort for every $R$ such that $R' \geq 0$ and $R'' \geq 0$ and $c\sigma + \mu$ is concave, then $F(x, \mu, \sigma)$ will satisfy condition (2.15).*

**Example 2.** *Analogously as in the previous example, we will have the following result. If a random variable $X$ that depends on a parameter $\theta$, with cdf $F(x, \theta)$, has the property that through a change of variable we have that $F(x, \theta) = F(r(\theta)x, 1)$, then, if $X$ depends on $e$ only through $\theta$ and $R(y \cdot r(\theta))/r(\theta)$ is convex and non-increasing in effort for every $R$ such that $R' \geq 0$ and $R'' \geq 0$, then*

- *$F(x, \theta)$ will satisfy condition (2.15).*

- *For any number of independent realizations of $X$, $n$, the distribution of the $k$th percentile as defined in (2.18), will satisfy condition (2.15).*

- *The "CVaR" non-Bernoulli utility function defined in section 2.2.2 will be concave in $e$ if (2.16) is satisfied.*

If conditions for the First Order Approach to be valid cannot be verified analitically, a heuristic to solve the optimization problem would then be to pose the problem assuming that FOA is valid, find the corresponding payment schedule $p^*(\cdot)$ and effort $e^*$ that maximize the providers utility and then verify the uniqueness of the stationary point in $e$ ex-post, given the payment schedule $p^*(\cdot)$. In fact, even if there is not a unique stationary point one needs only to verify that the value of $e^*$ is the global maximum of the utility of the provider, given $p^*(\cdot)$.

## 2.4.   SLA determination in a n-tier IT Service Scenario

SLA/contracts for IT Services often contain clauses regarding desired levels of response time. To provide a certain level of response time, a service provider has to use costly resources. The response time obtained as a result will still be stochastic around an average value. In our framework an optimal menu of contracts has to take into account the randomness of any particular measure of response time, and the characteristics of the different types of clients. Consider a context in which to meet a prescribed metric of Response Time a Service Provider has to provision computing resources, in the form of computer servers. This is frequently the case for Application Service Providers such as e-commerce sites). In this article, we look at single tiered services (e.g. Database services, web server utilities, Application server utilities etc.) The Response Time of such a single-tiered IT Service is modelled using a simple analytic queuing theory model. We suppose we have an IT Service that receives requests, that arrive according to a Poisson process of parameter $\lambda$. Servers handle the requests and their service time will also behave as a r.v. If at any given time all servers are occupied with requests, all requests that arrive thereafter will wait in queue to be serviced. If we suppose that the requests, that one server receives, behave as a Poisson of parameter $\tilde{\lambda}$ and the service time of the server is exponential of parameter $\mu$, it is a known result from queuing theory that the total Response Time will be exponentially distributed with parameter $\mu - \tilde{\lambda}$. If the workload is shared equally among the compute servers, then each server will receives requests at a rate $\lambda/e$, where $e$, the effort variable, is the number of servers, then the Response Time will distribute exponentially with parameter $\bar{\lambda}(e) = \mu - \frac{\lambda}{e}$ . This derivations is valid as long as $\mu > \lambda/e$.

SLA/contracts vary, in terms of the performance metric of response time that is used.

60

The quality variable that is appropriate to use in the developed framework will depend on the particular performance metric in which any particular contract is written. For example, if a contract specifies, that the hourly average Response Time is lower or equal than 25 ms, the quality variable we propose for this case would be the realized hourly Response Time. This quality variable will follow a probability distribution that can be derived under this framework in which each Response Time behaves exponentially. In general, a befitting quality variable will be one that appears explicitly in the SLA to specify a determined level of service, and in terms of which penalties will depend. In the previous example, the clause could specify penalties such as, the SP will pay a penalty of 10,000 if the average response time is between 25 and 35 ms and 20,000 if the average response time is between 35 and 45 ms.

In this work we focus on some likely SLA clauses: (a) contracts that specify a desired average of response time lower than $t$, over a period of length $\bar{T}$; (b) contracts that are in terms of percentiles, such as 95% of the requests have a response time lower than $t$, over a period of length $\bar{T}$ ; (c) combinations of the two previous types, such as 95% of the hourly averages have to be lower or equal than $t$, over a period of length $\bar{T}$.

## 2.4.1. Mean Response Time

If an SLA is written in terms of the Average Response Time, computed during a length of time $\bar{T}$. If the total number of requests is $n$ and the respective realized Response Times are $t_1, t_2, t_3, \ldots, t_n$, a possible Quality Variable, to fit this context, would be $-\bar{t}^n = -\frac{\sum_{i=1}^{n} t_i}{n}$, that is, minus the Average of Response Time [4]. The distribution of the average as an statistic depends on the original distribution of the sample.

### 2.4.1.1. Distribution of the Mean if Response Times are Exponentially Distributed

If we assume that the Response Times are exponentially distributed, the distribution of the mean Response Time, conditional on the total number of requests, $n$, will be a Gamma$(n, \frac{1}{n\lambda})$

$$f_{-\bar{t}^n}(t) = \frac{(\bar{\lambda}n)^n t^{n-1} e^{-\bar{\lambda}nt}}{\Gamma(n)} \tag{2.17}$$

---

[4]Note that utility has to be increasing in quality

If we know that the process of arrival of requests is also exponentially distributed, we can determine the distribution of the Average Response Time as

$$f_{-\bar{t}}(t) = \sum_{n=0}^{\infty} \frac{e^{-\lambda} e^{\bar{\lambda}nt} l^n \left(\bar{\lambda}n\right)^n t^{-1+n}}{\Gamma(n)\Gamma(1+n)}$$

However, for the sake of computational simplicity, we could estimate the quantity of requests by its mean, $\lambda\bar{T}$, in which case to determine the distribution of $\bar{t}$ we would have to replace $n = \lambda\bar{T}$ in (2.17).

This quality variable satisfies the conditions to use the FOA approach as it can be seen in Example 2.

### 2.4.1.2.  Normally Distributed Mean

If the sample is big enough, from the Central Limit Theorem, we can assume that the mean has a Normal Distribution. This is useful if we are not sure of the underlying distribution of each response time or of any other quality variable that we are analyzing. In this case, since each response, has mean $\mu = \frac{1}{\lambda(e)}$ and variance $\sigma^2 = \frac{1}{\lambda(e)^2}$, a normally distributed mean will be a $N(\frac{1}{\lambda}, \frac{1}{n\lambda^2})$. However, assuming normality of response time might not be so appropriate for this particular case because the Normal distribution takes on any value in the real line, and response time is positive. An alternative would be to use a truncated normal distribution, that is for some positive function $r(e)$ with $-r(e) + \mu \geq 0$ we will have that

$$f_{-\bar{t}}(t) = \begin{cases} \frac{e^{-\frac{(t+\mu)^2}{\sigma}}}{\frac{\sigma\sqrt{2\pi}}{d(e)}} & \text{if } t \in [-r(e) - \mu, r(e) - \mu] \\ 0 & \sim \end{cases}$$

Where $d(e) = \int_{-r(e)-\mu}^{r(e)-\mu} \frac{e^{\frac{(x+\mu)^2}{2\sigma^2}}}{\sigma\sqrt{2\pi}} = \int_{-\frac{r(e)}{\sigma}}^{\frac{r(e)}{\sigma}} \frac{e^{\frac{x^2}{2}}}{\sqrt{2\pi}}$. Note that if $r(e) = constant \cdot \sigma(e)$, $d(e) = d$ will not depend on $e$. In that case $-\bar{t}$ will satisfy conditions (2.15)-(2.16) as it is shown in the appendix, using Example 1.

## 2.4.2. Percentiles

It is common that a Service Level Requirement is in terms of the quantile of response time, such as "95% of the realizations of Response Time have to be less than 5 seconds". A quality variable for this type of contract could then be related to the realized 95 percentile of response time, measured during a time period.

We will use the following definition of the $k$th percentile, $P_k$. If $n$ is the number of requests, and $x_1, x_2, \ldots, x_n$ represent the ordered values of Response Time.

$$P_k = x_j \qquad \text{, where } j = round((n-1) \cdot \frac{k}{100} + 1) \qquad (2.18)$$

Let $T_k^n$ be the RV of the $k$th percentile given that there were $n$ requests during the measurement period. Let's assume that the distribution of Response Times are iid and if $F(t)$ is the cumulative distribution function of each Response Time. The $k$th percentile is an order statistic, therefore its distribution will be,

$$f_{T_k^n}(t) = \frac{d}{dt} \sum_{i=j}^{n} \mathbb{P}(t_1 \leq t, t_2 \leq t, \ldots, t_i \leq t, t_{i+1} \geq t, \ldots, t_n \geq t) = \frac{n!}{(j-1)!(n-j)!} F(t)^{j-1} (1-F(t))^{n-j} f(t).$$

$$(2.19)$$

where $j$ is given by (2.18).

### 2.4.2.1. Distribution of the Percentile if Response Times are Exponentially Distributed

If we assume that the Response Time of the servers system is exponentially distributed. In (2.19), we would have that $t_i \quad i \in \{1, 2, \ldots, n\}$ are r.v. iid, exponentially distributed, with parameter $\bar{\lambda} = \mu - \frac{\lambda}{e}$. This is not consistent with the fact that we are assuming that only $n$ events took place, however if $n$ is much bigger than the quantity of people in queue at any given time, the assumption becomes reasonable.

The probability distribution function of $T_k^n$ will be:

$$f_{T_k^n}(t) = \frac{n!}{(j-1)!(n-j)!} (1 - e^{-\bar{\lambda} \cdot t})^{j-1} (e^{-\bar{\lambda} \cdot t})^{n-j} \bar{\lambda} e^{-\bar{\lambda} \cdot t}.$$

where $j$ is defined by (2.18).

63

We will take $n$ to be the mean of arrivals, that is $n = \lambda \bar{T}$, where $\bar{T}$ is the length of the measurement period.

If we take the quality variable to be minus the $k$th percentile of response time, we would have that the utility of the client is increasing in quality. Under this conditions, from Example 2 it can be verified that the FOA conditions (2.15)-(2.16) will be fulfilled.

### 2.4.3. Contracts that are in terms of means and percentiles

Using what has been discussed above, we can easily derive the probability distributions for quality variables that have to be in terms of percentiles and means.
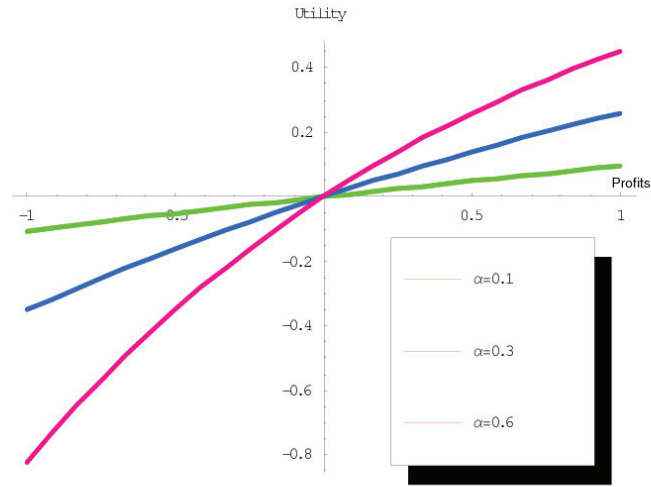
1. If a contract is in terms of percentiles of averages, such as 95% of the hourly averages have to be lower or equal than $t$, over a period of length $\bar{T}$, its probability distribution will be given by (2.19), where $F$ will be a Gamma or a Normal, depending on which distribution we choose to represent the mean. Later, using Examples 1 and 2, FOA conditions can be verified.

2. If a contract is in terms of averages of percentiles, such as the average of the hourly 95 percentiles have to be lower or equal than $t$, over a period of length $\bar{T}$, we know what the probability distribution of the percentiles will be and the distribution of the average of percentiles has to be determined. However, in this case, the Central Limit Theorem will tell us that if the average is taken over a big sample of percentiles, the average will be Normal, in which case we will need only to determine the mean and variance of the percentiles.

## 2.5. A Numerical Example

We computed optimal menus of contracts for different scenarios. The utility function that was used is known as the CARA (constant risk aversion) utility function. That is given a random variable $X$, that describes the behavior of uncertain profits, with probability density function $f_X(x)$ the utility of outcome $x$ will be $u(x) = (1 - e^{-\alpha x})$ the Expected Utility will then be $U(X) = \int u(x) f_X(x) dx$. The CARA Utility Function has constant risk aversion equal to $\alpha$. In figure 2.3, it can be seen that the greater the parameter $\alpha$ the more concave the CARA function is and the more the agent dislikes

risk. In our example we will have $U(p(g), e) = u(p(g) - \phi(e))$ where $\phi(e)$ is convex in $e$.

Figure 2.3: CARA Utility Function



## 2.5.1. Scenario: Mean Response Time and CARA Utility Function

We considered a case in which the quality variable is $q = -\bar{t}$, where $\bar{t}$ is the the realized Average Response Time, and is Gamma distributed. The $g$ function, which represents the monetary valuation that a client gives to the quality variable, was taken to be of the following form

$$g_{k,m,t}(x) = \begin{cases} m(x - t) + k; & \text{if } x \leq t \\ k & \text{if } x \geq t \end{cases}$$

Different types of clients are parameterized by having different values $k, m$ and $t$. The $g$ function will increase linearly in $-\bar{t}$, with slope $m$, until a point in which it becomes constant and equal to $k$. The reason to assume this form of $g$ function is that we assume that the clients value more quality up to a certain point in which they "saturate": greater quality does not increase his monetary utility any further. In this scenario we assume that the clients and the SP have CARA utility functions.

We will let the parameter of risk aversion $\alpha$ to vary for all agents. The value of $\bar{T}$ was taken to be 0.5 and we assume there is the same proportion of clients from each client type. The optimization problem was solved using the First Order Approach although conditions given in Lemma 1 were not verified. The validity of the approach was verified ex-post. The set of efforts per client will be bounded below. For practical reasons in the numerical computations, for each client type $\theta_j$, $e_{\theta_j}$ was taken such that $e_{\theta_j} \in [\frac{\lambda_{\theta_j}}{\mu_{\theta_j}} + \varepsilon, \infty)$, with $\varepsilon$ small. We don't include the multipliers that correspond to the lower bounds of effort in the First Order Condition that represents (2.10), because we assume that the solution of (2.10) will be interior, which is later verified in practice.

#### 2.5.1.1. Varying Risk Aversions

As a first analysis let's suppose that the SP is facing clients that have the same valuations of the service, that is, the same $g$ function, that is shown in figure 2.4, but they differ in the risk aversion parameters $\alpha$. Clients of type 1, 2 and 3 will have risk aversion equal to 3, 1.5 and 0.1, respectively. The values of $\bar{V}$ were varied with $\alpha$ keeping the certainty equivalent fixed. In tables 2.1 and 2.2 we present the parameters used for the computation of the optimal menu of contracts and in table 2.3 we present the optimal menu of contracts obtained. Parameters $\mu$ and $\lambda$ didn't change across client types and were 5 and 100 respectively.

Figure 2.4: Example of "g function" (m=10, k=150000, t=-2.5)
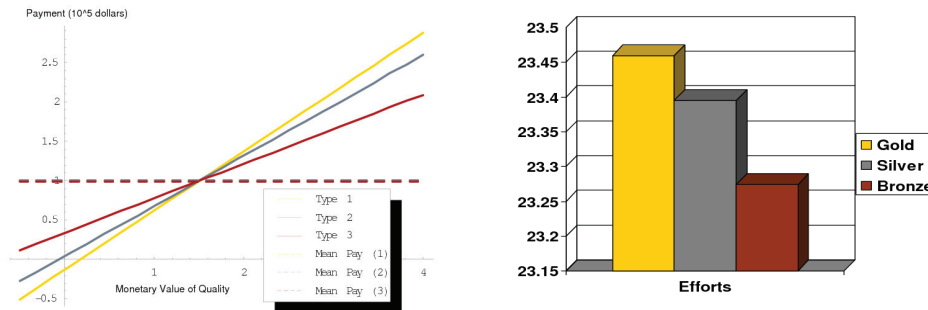
Table 2.1: Parameters Service Provider

| Cost | Parameter of SP |
|------|-----------------|
| 2200 | 1.5 |

Table 2.2: Parameters Clients

| Type of Client | $\bar{V}$ | Risk Parameter of Clients | | |
|----------------|-----------|---------------------------|--|--|
| Type 1 | 0.78 | 3 | | |
| Type 2 | 0.53 | 1.5 | | |
| Type 3 | 0.049 | 0.1 | | |

In figure 2.5 we present a plot of the contracts obtained. The dashed lines represent the mean payments that each client type will make. The clients of type 1 who are the most risk averse will pay more for qualities that are above the mean than the other client types, and will pay less for qualities under the mean. That is,the more averse clients are better insured to variatons of quality. The opposite is true for clients of type 3 who are the least risk averse.

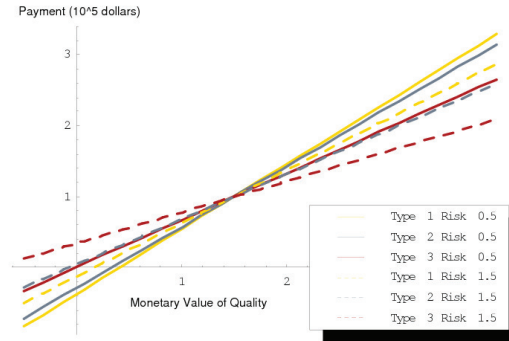Figure 2.5: Different Risk Aversions



If now we change the risk aversion parameter of the service provider from 1.5 to 0.5, the slopes of the linear contracts become higher for all three client types, as it is shown in figure 2.6. It becomes more costly to generate the right incentives to the SP when she is less risk averse, and in order to induce effort penalties must be more stringent.

In order to assess the optimality of the menu of contracts presented we pose ourselves the question of what the profits would be if the SP offered a different menu of contracts. If the SP offers only one contract to all clients, the self-selection constraints will be

Table 2.3: Optimal Menu of Linear Contracts

| Type of Client | Level of Effort | Slope | Intercept |
|:---:|:---:|:---:|:---:|
| **Type 1** | 23.46 | 0.75 | -0.13 |
| **Type 2** | 23.39 | 0.64 | 0.04 |
| **Type 3** | 23.27 | 0.44 | 0.34 |

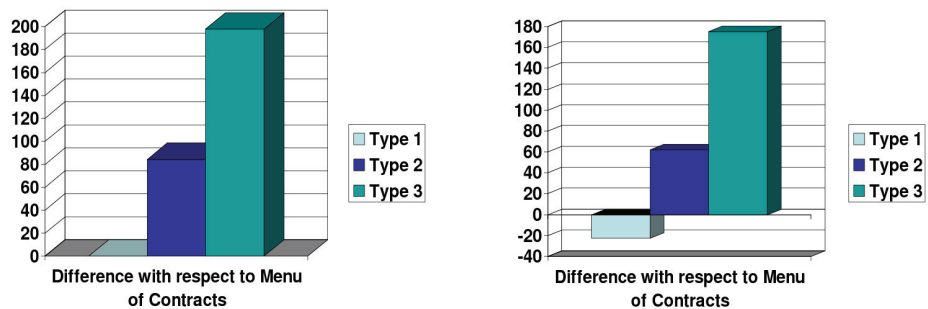Figure 2.6: Changing Risk Parameter of Service Provider



satisfied trivially. If we were to offer only one of the contracts of the three presented in table 2.3 it would have to be the contract offered to the clients of type 1, because the contract intended for their own type is the only contract that clients of type 1 are willing to accept. In figure 2.7 (left) we present the losses that the SP would experience if she offered such contract, with respect to the optimal menu of contracts in table 2.3. Profits made from payments of clients of type 1 will remain constant, while profits from the other two client types will be lower. In table 2.4 we present the optimal contract that the SP would offer to each client type if she could know which type is which, we refer to this as the "perfect discrimination case". If the SP were to offer one of the three, as before, the only contract that would be accepted by type 1 clients would be the one that is optimal for their client type. In figure 2.7 (right) we present the differences between this contract and the optimal menu of contracts in table 2.3. The losses are lower in this case, but they are positive. Note also that these client types are identical in every respect except for their risk aversion factors.

Table 2.4: Optimal Contracts with Perfect Discrimination

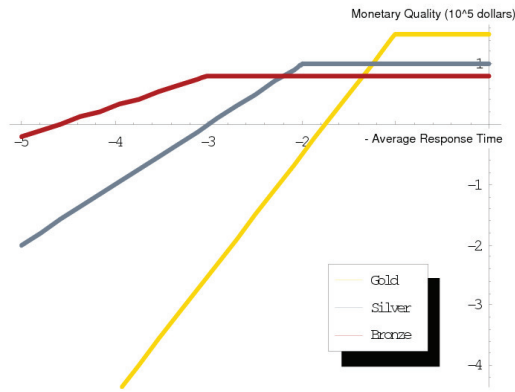| Type of Client | Level of Effort | Slope | Intercept | Profits per Client Type |
|---|---|---|---|---|
| Type 1 | 23.44 | 0.70 | -0.05 | 48350.25 |
| Type 2 | 23.37 | 0.59 | 0.12 | 48443.54 |
| Type 3 | 23.28 | 0.44 | 0.34 | 48573.78 |

Figure 2.7: Comparisons with Optimal Menu of Contracts

### 2.5.1.2. Gold, Bronze and Silver Clients

The framework presented above allows us to compute optimal menu of contracts for clients who were different in many dimensions. In figure 2.8 we present the $g$ function for three hypothetical types of clients, the first type will have parameters $m = 1, k = 1, t = -2$, the second, $m = 2, k = \frac{3}{2}, t = -1$, and the third, $m = \frac{1}{2}, k = 0.8, t = -3$. We will refer to them as type "Gold", "Silver" and "Bronze", respectively, and we will assume that there is the same amount of clients of each type.

The second type values high quality of service more than the other two (higher $k$), however, his profits decrease faster as quality goes down, also his saturation point is higher. The third type requires lower levels of service, he has a low saturations point, and his profits don't decrease very fast if quality becomes lower, he also values the highest quality less. The first one would be the "middle" type. We will give each client type a different parameter $\alpha$ and a different value of $\bar{V}$.

Figure 2.8: Gold, Silver and Bronze Clients



In tables 2.5 and 2.6 we present the parameters for the SP and for each client type. The Gold clients will have a higher risk aversion parameter and $\bar{V}$, and the Bronze clients will have the lowest value for those two parameters.

In figure 2.9 we present the optimal linear contract and the respective levels of effort and in dashed lines, the mean payments. The SP will make more more profits from clients of type Gold, as it can be seen in figure 2.10.
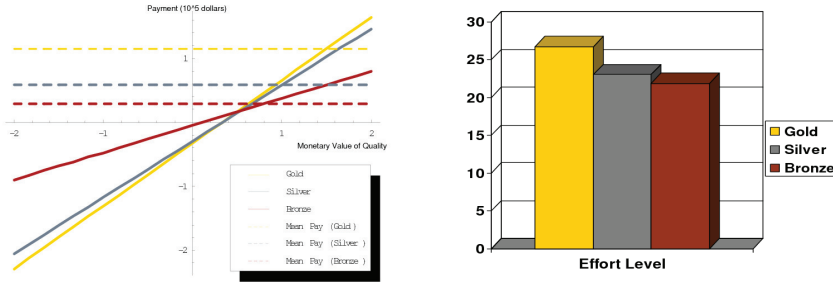
Table 2.5: Parameters Service Provider

| Cost | Parameter of SP |
|------|-----------------|
| 1000 | 0.1 |

Table 2.6: Parameters Clients

| Type of Client | $\bar{V}$ | mu | lambda | Risk Parameter of Clients |
|----------------|-----------|-----|--------|---------------------------|
| **Gold** | 0.1 | 5 | 100 | 0.3 |
| **Silver** | 0.08 | 5 | 100 | 0.2 |
| **Bronze** | 0.05 | 5 | 100 | 0.1 |

Figure 2.9: Optimal Contract 1



The contracts have to be translated into (true) quality. In figure 2.11 we present the contracts in terms of (true) quality, and in dashed lines the mean payments for each client type. Note that the mean payments will be very close to the highest payments possible. This is because the efforts assigned to each client, for this example, will deliver a quality inside the area of saturation with high probability.

## 2.5.2. Scenario: Mean Response Time and Exp-Var Function

Very similar results are obtained when the Exp-Var utility function is used. For this function we didn't verify sufficient conditions for the FOA approach, but its validity is confirmed ex-post. The value of $\bar{T}$ was taken to be 2 and we assume there is the same proportion of clients from each client type. The "g" function that was used is the same as in the first part of the previous scenario, and is showed in figure 2.4. In tables

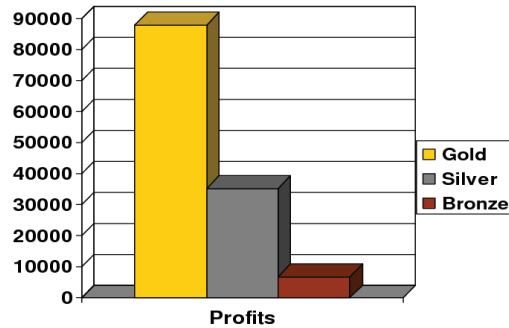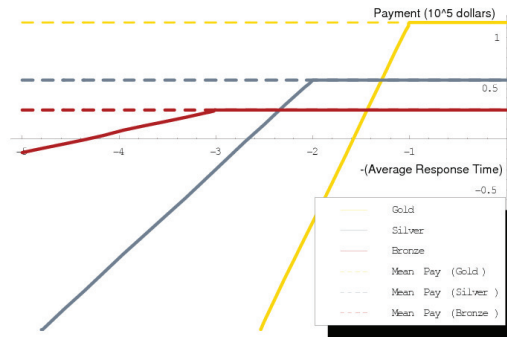Figure 2.10: Profits from each client type



Figure 2.11: Quality vs Payment



2.7 and 2.8 and in figure 2.12 we represent graphical representations of the contracts obtained.

Table 2.7: Parameters Service Provider

| Cost | Parameter of SP |
|------|-----------------|
| 2200 | 6 |

Note that it is not straightforward to make comparisons between the contracts obtained using the CARA utility function and the Exp-Var. If the agents have different utility functions, their preferences will be different, even if both utility functions were appropriate to represent their preferences, a calibration of the parameters has to be made. However, we can see from the results obtained (see table 2.9 and figure 2.12) that we

Table 2.8: Parameters Clients

| Type of Client | $\bar{V}$ | mu | lambda | Risk Parameter of Clients |
|:---:|:---:|:---:|:---:|:---:|
| Type 1 | 0.5 | 5 | 100 | 3 |
| Type 2 | 0.5 | 5 | 100 | 1.5 |
| Type 3 | 0.5 | 5 | 100 | 0.1 |

Table 2.9: Optimal Contracts

| Type of Client | Level of Effort | Slope | Intercept |
|:---:|:---:|:---:|:---:|
| Type 1 | 23.02 | 5.95 | -7.93 |
| Type 2 | 22.99 | 4.78 | -6.17 |
| Type 3 | 22.89 | 1.66 | -1.49 |

get analogous interpretations for the contracts for different clients given their risk parameter.

If now we decrease the risk parameter of SP from 6 to 2, we can see in figure 2.13 that the slopes of the contracts for every client will increase, this has a similar interpretation as in the previous case.
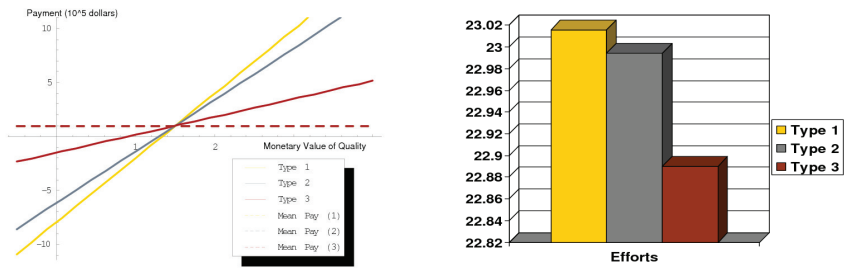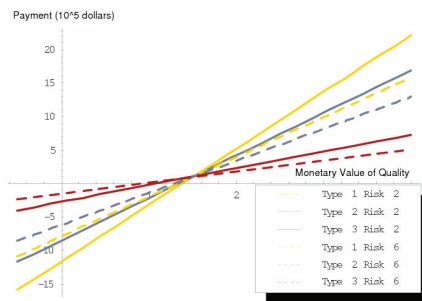
Figure 2.12: Optimal Contract 2



Figure 2.13: Optimal Contract 3

## 2.6.  Conclusions

In this paper we have applied contract theory to find an optimal SLA/Contract given the characteristics of the client and service provider. We also extended the basic model to take into account the possibility that the service provider will offer different contracts to different client types,in order to cover varied necessities and to extract more profits. We demonstrated this through a model of single-tiered IT services. Through the numerical example,in which a number of insights, that were consistent with economical intuition, were developed. We also analyzed the conditions under which a First order Approach can be used, based on the literature of the principal-agent problem.   In practice, the usefulness of our model, if calibrated correctly, is to give benchmarks for future contracts, in each stage of an eventual negotiation process.  For calibrations purposes, information of past contracts can be used. It is important to note that there still are aspects of the determination of SLA/Contracts that are not tackled here. For instance, our framework requires that the agents have a great deal information about each other.  This is not generally the case, since, in reality, there might exist a gap of information between the agents.  A possible line for future research that could be accounted for, and in turn, this could also give us a greater understanding on how the bargaining (negotiation) process takes place.

# Bibliography

K. Appleby, S. Fakhouri, L. Fong, G. Goldszmidt, M. Kalantar, S. Krishnakumar, D.P. Pazel, J. Pershing, and B. Rochwerger. 2001. SLA Based Management of a Computing Utility. *Integrated Network Management Proceedings.*

M.J. Buco, R.N. Chang, L.Z. Luan, C. Ward, J.L. Wolf, P.S. Yu, 2004, *IBM Systems Journal, VOL 43, NO 1, 2004*, 159–178.

George Candea and Armando Fox. 2002. A Utility-Centered Approach to Building Dependable Infrastructure Services. *Proceedings of the 10th ACM SIGOPS European Workshop (EW-2002).* 213–218, Saint-Émilion, France, September 2002.

P. Chone and Jean-Charles Rochet. 1998. Ironing, Sweeping and Multidimensional Screening. *Econometrica*, vol. 66, n. 4, 1998, 783–826.

Haluk Demirkan, Michael Goul, Daniel S. Soper, 2005, Service Level Agreement Negotiation: A Theory-based Exploratory Study as a Starting Point for Identifying Negotiation Support System Requirements, *Proceedings of the 38th Hawaii International Conference on System Sciences - 2005.*

Jewitt, Ian, 1988, Justifying the First-Order Approach to Principal-Agent Problems, *Econometrica 56* ,(5), 1117–1190.

R. Tyrrell Rockafellar, Stanislav Uryasev, Michael Zabarankin, 2002, Deviation Measures In Risk Analysis And Optimization, *RESEARCH REPORT # 2002-7* Risk Management and Financial Engineering Lab Center for Applied Optimization Department of Industrial and Systems Engineering University of Florida, Gainesville.

Akhil Sahai, Anna Durante, Vijay Machiraju, 2002, *Software Technology Laboratory, HP Laboratories*, Palo Alto, HPL-2001-310 (R.1).

Stole, Lars, 2001, Lectures on the Theory of Contracts and Organizations.