#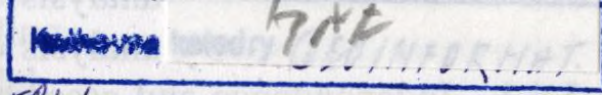 Contents