

1 **Abstract**

2 **Background:** Researchers interested in the effect of health on various life outcomes often use self-
3 reported health and disease as an indicator of true, underlying health status. However, the validity
4 of reporting is questionable as it relies on the awareness, recall bias and social desirability.
5 Accordingly, the measured biomarker is generally regarded as a more precise indication of the
6 disease.

7 **Objectives:** The study aimed to examine the discrepancy between the reporting and biomarkers
8 of hypertension and diabetes in the contemporary China, and explore sociodemographic
9 characteristics that are correlated with misreporting.

10 **Methods:** Using data from the third wave of China Health and Retirement Longitudinal Study
11 (CHARLS), we selected individuals aged 40-85 years old who participated in both a health
12 interview survey and a biomarker examination. Sensitivity, specificity, false negative reporting
13 and false positive reporting were used as measurements of (dis)agreements or (in)validity. Binary
14 and multinomial logistic regression were used to estimate under-report or over-report of
15 hypertension and diabetes.

16 **Results:** Self-reported hypertension and diabetes showed low sensitivity (71.98% and 49.21%,
17 respectively) but high specificity (93.71% and 98.05%, respectively). False positive reporting of
18 hypertension and diabetes were 3.85% and 1.67%, while false negative reports were extremely
19 high at 10.85% and 7.38%. Education degree, hukou, age and gender affected both the specific
20 error and the overall error of reporting hypertension and diabetes, but there were some differences
21 in the magnitude and direction.

22 **Conclusion:** Self-reported conditions underestimate the disease burden of hypertension and
23 diabetes. Adding objective measurements into social survey could improve data accuracy allowing
[Insert Running title of <72 characters]

1 better understanding of socioeconomic inequalities in health, especially collecting biological
2 indicators for populations with limited access to regular healthcare in China. Furthermore, there is
3 an urgent need to provide basic health education and physical examination to citizens, to facilitate
4 access to healthcare and make focused interventions to lower the incidence and unawareness of
5 disease in China.

6 **Keywords:** Self-reporting, Biomarker, Discrepancy, Sociodemographic characteristics

8 **Introduction**

9 Hypertension and diabetes are two well-known risk factors of cardiovascular disease, the
10 leading cause of death worldwide with 17.79 million deaths in 2017(Ritchie and Roser 2019).
11 Prevalence estimates of high-quality based on biomedical measurements of both diseases are
12 needed for monitoring cardiovascular disease risks and planning public health interventions and
13 prevention. Due to the high cost and long-time collection of biomedical data, economists and
14 demographers have relied heavily on self-reported disease of hypertension and diabetes to estimate
15 the prevalence and disease burden. However, recent research has raised doubts about the reporting
16 error of self-reported disease. Such error can occur for a number of reasons. For example, survey
17 respondents may report their disease differently depending on last specific behavior, socially
18 driven conceptions of what ‘chronic disease’ means, expectations of their own health, using of
19 healthcare, and comprehension of the actual survey questions (Murray and Chen 1992, Newell,
20 Girgis et al. 1999).

21 Many a study has attempted to assess the value of a self-reported disease by comparing self-
22 reports with objective assessment, which has many advantages: precision of measurement, reliable,
23 and less biased than questionnaires (Mayeux, 2004). The criterion shows the extent between self-
[Insert Running title of <72 characters]

1 reports and biomarker and is usually measured by sensitivity and specificity of the errors in a given
2 group, as well as false reporting of the overall error. Sensitivity was defined as the percentage of
3 respondents who reported having hypertension / diabetes among those with biomedical
4 hypertension / diabetes. This value was thus equivalent to ‘hypertension / diabetes awareness
5 among those with diseases; Specificity was defined as the percentage of individuals who reported
6 not having hypertension / diabetes among those with ‘normal’ or ‘healthy’ biomedical
7 measurements; False negative reporting was defined as who reported not having hypertension /
8 diabetes but were diagnosed hypertension / diabetes compared with those people of correct
9 reporting, and false positive reporting was defined as who self-reported having hypertension /
10 diabetes but were not diagnosed hypertension / diabetes compared with those people of correct
11 reporting. Evidence from developed countries indicated that there was a big gap between self-
12 reported disease and biomedical disease of hypertension and diabetes and they may differ between
13 socioeconomic groups. During 2011-2014, nearly 16% of adults with hypertension were unaware
14 of their hypertension, and a higher percentage of non-Hispanic Asian (24.7%) and Hispanic
15 (20.2%) adults than non-Hispanic white (14.9%) and non-Hispanic black (14.7%) adults were
16 unaware of their hypertension status in U.S (Paulose-Ram, Gu et al. 2017). The National Diabetes
17 Statistics Report showed 23.8% of American adults with diabetes were undiagnosed and
18 prevalence varied significantly by educational level (Control and Prevention 2017). Baker et al.
19 compared the self-reported status from a Canadian household survey with health administrative
20 data in the province of Ontario and they found false negative rates were over 50% (Baker, Stabile
21 et al. 2004). Johnston et al. used the health survey from England to examine the gap between self-
22 reported and clinically tested hypertension and found 28% under-reporting on average, and
23 estimated an attenuation bias of 68% (Johnston and Propper 2009). More educated and higher
[Insert Running title of <72 characters]

1 socioeconomic individuals may have a better understanding of health information and were more
2 capable to answer survey questions on disease diagnosis (Kislaya, Tolonen et al. 2019).

3 Although developed countries have witnessed substantial disagreements between self-reports
4 and objective diagnosis of the two diseases and found it may differ between socioeconomic groups,
5 the performance of developing countries were unclear due to the data limitation. A handful of
6 studies have investigated the incidence of disease awareness in local areas. Onur et al. used
7 Longitudinal Aging Study in India (LASI) data of four selected states and found the average under-
8 reporting rate of hypertension was 26%, and only 4% reported having lung disease while 43% of
9 the sample tested positive (Onur and Velamuri 2018). Most patients with hypertension and other
10 NCDs were unaware of their condition, and hypertension in treated patients was mostly
11 uncontrolled (Kavishe, Biraro et al. 2015). Huang et al. (2017) looked into the prevalence,
12 awareness, treatment, and control of hypertension among very elderly in southwestern China and
13 they found the prevalence of hypertension is predominantly high, whereas awareness, treatment,
14 and control rates are considerably low (Huang, Xu et al. 2017). There is still debate regarding the
15 consistency and direction of the social gradient and health in developing countries (Brasher et al.
16 2017), which means socioeconomic status might differ in developing countries compared to other
17 contexts. One possible explanation is that economic development takes place unequally across
18 regions, and so does progress in the awareness of a disease. In addition, different population vary
19 in exposure to health enhancing or health-damaging factors. For example, the household
20 registration (hukou) system required all Chinese households to be registered in the locale where
21 they resided and categorized as either agricultural or nonagricultural (rural or urban) status (Wu,
22 2019), which divided China into two separated societies, with the majority of the population
23 confined in the rural and entitled to few of the rights and benefits that the socialist state conferred
[Insert Running title of <72 characters]

1 on urban residents, thus creating not only a spatial stratification between the countryside and the
2 cities but also two unequal classes of Chinese citizens (Solinger 1999, Wu & Treiman 2004).

3 With approximately 20% of the world's population, China is experiencing rapid population
4 aging and an epidemiological transition moving from the primacy of acute infectious and
5 deficiency diseases to the increasing dominance of non-communicable and chronic conditions
6 (Brasher, George et al. 2017). Report on Nutrition and Chronic Diseases of Chinese Residents
7 (2015) shows that the prevalence rate of hypertension and diabetes in China has reached 25.2%
8 and 9.7% over the age of 18, respectively. The lack of medical insurance and the high cost of the
9 health care market may lead to incorrect perception of disease prevalence, which means under-
10 diagnosis of chronic disease is not uncommon and self-reported diseases are inaccurate in China.
11 However, whether and to what extent the discrepancy between self-reported and biological
12 hypertension and diabetes is not clear in China. Meanwhile, diabetes mellitus caused the most
13 death from the 19th in 1990 rises to 8th in 2017 and high systolic blood pressure was the leading
14 risk factors contributing to deaths and DALYs in China (Zhou, Wang et al. 2019). Therefore,
15 understanding the real situation of the two diseases, estimating the burden of diseases and
16 intervening the risk of death from diseases are very important in China.

17 In this paper, we compare the discrepancy between self-reported and biomarker of
18 hypertension and diabetes for the same individuals, using a nationally representative database on
19 older individuals in China. Our paper also analyzes the sociodemographic characteristics of miss
20 reporting in all and specific groups. To my knowledge, the present study is the first to address this
21 issue by offering four indicators' comparison on self-reported disease and biological outcomes

[Insert Running title of <72 characters]

1 and its sociodemographic characteristics difference measured on the representative samples in
2 China.

3

4 **Data and methods**

5 **Data**

6 Data for this analysis are from China Health and Retirement Longitudinal Study (CHARLS),
7 which is a nationally representative survey of the middle aged and elderly in China, designed by
8 the National School for Development (China Center for Economic Research) together with the
9 Institute for Social Science Survey at Peking University. The baseline national wave of CHARLS
10 was being fielded in 2011 and included about 10,000 households and 17,500 individuals in 150
11 counties/districts and 450 communities. The multistage sample was drawn at each stage based on
12 probability-proportional-to-size random-sampling procedures. The survey collected detailed
13 demographic backgrounds, socioeconomic information and health status and functioning. The
14 analysis draws on data from Wave 3 (2015), because it collects reliable venous blood samples
15 which allow us to compare the discrepancy between self-reported health and underlying biomarker
16 levels among CHARLS respondents. Detailed information about the CHARLS blood sampling
17 procedure and data quality management has been published previously (Chen, Crimmins et al.
18 2019).

19 We restricted the sample as people 40 to 85 years old with valid self-reported diseases and
20 biomarker values. The Wave 3 (2015) included 20,284 respondents asked to consent to a venous
21 blood draw and blood pressure measurement; 13,013 provided venous blood and 16,406 provided
22 blood pressure information. Combined with those people who also provided detailed educational

[Insert Running title of <72 characters]

1 information, the final sample size about hypertension and diabetes were 14,457 and 12,189,
2 respectively. Characteristics of the sample are summarized in Table 1^①.

3 **Measures**

4 *Self-reported and biomedical measures of hypertension and diabetes*

5 Self-reported data on hypertension and diabetes were obtained by the question, ‘Have you
6 been diagnosed with hypertension / diabetes by a doctor?’. Each respondent who answered no,
7 was continued to ask “do you know whether you have hypertension by yourself”^②. If a respondent
8 answered in the affirmative for any of two questions, we defined the self-reported hypertension or
9 diabetes as 1, otherwise as 0. Biomedical blood pressure was measured three times (approximately
10 45 s apart) on a single occasion, using an electronic monitor. The average of these blood pressure
11 readings was used to determine each respondent’s blood pressure level. Hypertension was defined
12 as a systolic blood pressure ≥ 140 mm Hg and/or a diastolic blood pressure ≥ 90 mm Hg and/or
13 current use of antihypertensive medication, following the WHO guideline (Alwan 2011);
14 Biomedical diabetes was measured by venous blood data which provided glycated hemoglobin
15 (HbA1c). The diagnostic criterion for diabetes in our study was defined as HbA1c values $\geq 6.5\%$.
16 If a respondent’s glycated hemoglobin was over 6.5%, we defined the biomedical diabetes as 1,
17 otherwise as 0. Although HbA1c may not be the most widely used screening test, it has been

① The distribution was very similar and no bias was found.

② The question was only asked about samples that had not been diagnosed with hypertension by a doctor.

[Insert Running title of <72 characters]

1 suggested as an alternative means of screening for diabetes and has been used in this way in many
2 surveys (Bennett, Guo et al. 2007).

3 *Economic resources*

4 Educational attainment was measured by three levels. This variable indicated the highest
5 education degree attained by respondents in the survey point. The lowest category comprised
6 individuals holding no formal education (illiterate); intermediate education ranged from not
7 finishing primary school, home school to elementary school; and the highest category included
8 respondents holding middle school or above. Hukou, as one of the most important redistributive
9 institutions under Chinese state socialism, was a time-varying covariate that was measured whether
10 the respondent held a rural hukou (0 = No, 1 = Yes).

11 *Health behavior*

12 Drinking was a 5-category variable indicating the frequency of drinking last year: none
13 (coded 1), less than 3 days/month (coded 2), less than 3 days/week (coded 3), 4 to 6 days/week or
14 daily (coded 4), twice a day or above (coded 5). Smoking was a continuous variable indicating the
15 frequency of cigarettes/day, which ranged from 0 to 100.

16 *Demographic characteristics*

17 Gender was a binary variable: male (coded 1), female (coded 0). Age was a 5-category
18 variable ranging from 40 to 49(coded 1), 50-59 (coded 2), 60 to 69(coded 3), 70-79 (coded 4), 80
19 and above (coded 5). Marital status was a time-varying covariate indicating whether the respondent

[Insert Running title of <72 characters]

1 was in a marriage status: separated, divorced, widowed and never married (coded 0), married or
2 partnered (coded 1).

3 **Analytic strategy**

4 Our first step was to assess the difference in prevalence estimates based on two data collection
5 methods, the prevalence of hypertension and diabetes were calculated according to self-reported
6 information, as well as according to the results of biomedical measurements obtained from the
7 CHARLS. To assess the accuracy of self-reported data, sensitivity, specificity, false negative
8 reporting and false positive reporting were calculated, respectively.

9 Both only sensitivity and specificity were of no practical use when it came to helping the
10 clinician estimate the probability of disease in individual patients (Akobeng, Ramanan et al. 2006).
11 In addition, sensitivity and specificity assessed only individual errors in diagnosed or undiagnosed
12 diseases, respectively, but not overall errors. We identified both total error and the specific error
13 and assessed sociodemographic characteristics that are correlated with misreporting (sensitivity,
14 specificity, false negative reporting and false positive reporting), and binary and multinomial
15 logistic regression analysis were considered, controlling for education, hukou, drinking, smoking
16 age, gender and marital status. As the total error outcome (correct reporting, false negative
17 reporting and false positive reporting) has more than two categories. The model has 2 equations as
18 follows[®]:

[®] In order to assess the differences in communities and provinces, we also adopted random intercept model, the results showed that variations between communities and provinces were not statistically significant, so this article only presented the result of the common model.

[Insert Running title of <72 characters]

$$1 \quad \ln\left(\frac{\Pr(Y_i = 2)}{\Pr(Y_i = 1)}\right) = \beta_0 + \beta_1 \text{Economics} + \beta_2 \text{Behaviors} + \beta_3 \text{Demographics} \dots + \varepsilon$$

$$2 \quad \ln\left(\frac{\Pr(Y_i = 3)}{\Pr(Y_i = 1)}\right) = \beta_0 + \beta_1 \text{Economics} + \beta_2 \text{Behaviors} + \beta_3 \text{Demographics} \dots + \varepsilon$$

3

4 **Results**

5 **Sensitivity, specificity and false reporting of hypertension and diabetes**

6 The prevalence of hypertension was 38.71% according to biomedical test and 31.72%
 7 according to the self-reported data, indicating that self-reporting led to an underestimation of
 8 hypertension by 18.06%. Likewise, the prevalence of diabetes was 14.52% according to the
 9 biomedical data and 8.81% according to self-reports, indicating an underestimated prevalence of
 10 diabetes by 39.3% from self-reported data.

11 Both the prevalence of self-reporting and objective hypertension and diabetes increased over
 12 age in China, and biomedical hypertension and diabetes rose considerably faster with age than
 13 self-reporting, which means the gap between self-reporting and objective hypertension increases
 14 with age as shown in Figure 1 and 2. This is suggestive of undiagnosed hypertension and diabetes
 15 becoming more of a problem with individuals' age. Since China has the greatest number of oldest
 16 old adults in the world according to United Nations data, and undiagnosed high blood pressure and
 17 diabetes may become common over time.

18 Table 2 provides the sensitivity, specificity and false reporting of self-reported hypertension
 19 and diabetes compared with biomedical data. The overall sensitivity and specificity of self-
 20 reported hypertension were 71.98% and 93.71%, which meant 28.02% of people didn't know they
 21 had hypertension, and 6.31% of people thought they had hypertension. The false reporting of

[Insert Running title of <72 characters]

1 hypertension was 14.7%, specifically false positive reporting of hypertension was 3.85% and false
2 negative reporting was 10.85%. For diabetes, the overall sensitivity and specificity were 49.21 %
3 and 98.05%, which meant over 50% of people didn't know they had hypertension, and only less
4 than 2% of people thought they had diabetes. The false reporting of diabetes was 9.05%,
5 specifically false positive reporting of hypertension was 1.67% and false negative reporting was
6 7.38%, which meant less than 10% of people misreporting diabetes status.

7 Comparing the four indicators above of hypertension and diabetes, we found that the overall
8 miss reporting is different from the group error. Taken together, these results were suggestive of a
9 substantial public health problem of undiagnosed hypertension and diabetes in China. Note that
10 our use of self-reported hypertension and diabetes will increase the degree of underdiagnosis and
11 underestimate true disease burden in the population substantially.

12 **Sociodemographic characteristics of sensitivity, specificity and false reporting**

13 The estimated effects on our two sets of the discrepancy between self-reported hypertension
14 / diabetes and underlying biomarker using binary and multinomial logistic regression is shown in
15 Table 3 and 4. Columns 1 and 2 show the estimated effects of predictor variables on the outcome,
16 sensitivity and specificity between self-reported health and biomarker using binary logistic
17 regression. Columns 4 and 5-6 show the estimated binary and multinomial effects of predictor
18 variables on the reporting, but with reporting divided into three variables: correct reporting, false
19 negative and positive reporting. We transformed the coefficients into relative odds ratios by
20 exponentiation of the original coefficients.

21 First, we look into the reporting error for a particular group. Respondent characteristics
22 associated with sensitivity found that individuals with higher educational degree, urban hukou,
23 aged ≥ 50 years, female, and having healthy lifestyle, were strongly and independently associated
[Insert Running title of <72 characters]

1 with having more accurately self-reported hypertension than their counterparts among those with
2 real hypertension. The results in Table 3 suggest that no respondent characteristic was significantly
3 associated with more accurate reporting in specificity except for age group and healthy lifestyle,
4 where elderly people were more likely to erroneously report the absence of hypertension than those
5 younger than 60 years of age, and people who drank alcohol every day were more likely to report
6 errors among those without hypertension. However, smoking may more likely to correctly report
7 the absence of hypertension.

8 Next, we will identify the overall error. The likelihood of reporting errors decreases with
9 education, but not significantly. Urban hukou, no drinking, female, younger age and married
10 people were strongly and independently associated with correct reporting (Columns 3 of Table 3).
11 The coefficient for rural hukou was 0.877, meaning that compared with urban hukou, rural hukou
12 was a 12.3% increase in miss reporting. Compared with non-drinkers, people who drank every day
13 or more were more likely to report errors. Men were more likely to erroneously report than women.
14 Elderly people were more easily to erroneously reporting than those younger than 50 years of age,
15 and the error rate increased with age. Unmarried people were more prone to miss reporting than
16 those married people. Specifically, educational level had a significant effect on the risks of false
17 negative reporting but not significantly on false positive reporting (Columns 4–5 of Table 3). The
18 propensity to false negative reporting went down significantly with educational gradient. The
19 coefficients for primary education and secondary education and above were 0.875 and 0.829,
20 respectively, meaning that compared with the illiterate, primary education had a 12.5% decrease
21 in the rate of false negative reporting while secondary education and above had a 17.1% decrease.
22 Although people with higher educational attainment reported higher false positive than their
23 counterparts, the effect was not significant. Compared with the illiterate, people with primary
[Insert Running title of <72 characters]

1 education had a 11.6% increase in the rate of false positive reporting while having secondary
2 education had a 14.8% increase. Individuals, having a rural hukou, aged ≥ 50 years, male, and
3 having unhealthy lifestyle, are strongly associated with false negative reporting. False positive
4 reports were almost not related to sociodemographic characteristics.

5 For self-reported diabetes, education, hukou, drinking and age were associated factors with
6 sensitivity. Aged participants with higher levels of education, having an urban hukou and small
7 drinking were also more likely to accurately self-report diabetes than were their respective
8 counterparts among those with real diabetes. Multivariate analyses showed that younger people
9 with a rural hukou and small drinking had slightly more accurate reporting the absence of diabetes
10 than their counterpart among those without diabetes.

11 Our indicator of education had almost no statistically significant effect on false reporting of
12 diabetes, except that secondary education and above might slightly reduce false negative reporting
13 (Columns 3-5 of Table 4). Contrary to our expectations, false reporting between self-reported
14 diabetes and biomedical diabetes did not depend on educational level, which was similar to
15 previous research (Al Shamsi and Almutairi 2018). Educational level didn't affect the
16 disagreement of measurements for diabetes, although educational attainment was associated with
17 the prevalence of diabetes. Compared with people of urban hukou, those of rural hukou had a 18.6%
18 increase in the rate of correct reporting due to the lower prevalence of diabetes in rural area. Elderly
19 people were more prone to false reporting and false negative reporting than those younger than 50
20 years of age, and the rate went up significantly with age. To be specific, 50 to 59 years old people
21 had a 37.1% increase in false negative reporting while 60 to 69 years old had a 67.6% increase, 70
22 to 79 years old had an 96.7% increase and 80 and above years old had a 122.5% increase in false
23 negative reporting rate.

[Insert Running title of <72 characters]

1 The results confirmed that education degree, hukou, age and gender affect both the specific
2 error and the overall error of reporting hypertension and diabetes, but there are some differences
3 in the magnitude and direction.

4 **Discussion**

5 Using data from China adults aged between 40 and above, we analyzed two diseases that
6 were commonly used in clinical evaluations of health-related risk, one of which was fitness
7 biomarker while the other was disease risk biomarker, and compared the discrepancy between self-
8 reports and biomedical measurements. We found a large difference in the percentage of the sample
9 who reported having hypertension and diabetes (31.72% and 8.81%) relative to those who were
10 measured to have two diseases (38.71% and 14.52%). As the published report from China say,
11 solely on self-reported measures of disease will tend to underestimate the true extent of the disease
12 burden in contemporary China, and we do find self-reporting led to an underestimation of
13 hypertension and diabetes by 18.06% and 39.03%, respectively. Due to the more complex
14 collection methods and higher cost, we find that disease risk biomarker (diabetes) is more prone
15 to underestimation, which is really not conducive to disease assessment and intervention.

16 We also examined the sociodemographic characteristics that are correlated with misreporting
17 by using four indexes (sensitivity, specificity, false negative reporting and false positive reporting).
18 For hypertension, we find that education degree, hukou, age and gender affect both the specific
19 error and the overall error of reporting hypertension, but there are some differences in the
20 magnitude and direction. Educational attainment was an important explanatory factor, and had a
21 significant impact on sensitivity and false negative reporting, with each additional gradient
22 decreased in education reducing the probability of committing false negative report of
23 hypertension by about 12.5 and 17.1 percentage points, respectively, which meant the least
[Insert Running title of <72 characters]

1 educated people were tend to underestimate their disease burden than their more educated peers.
2 Besides, Hukou, as one of the most important redistributive institutions under Chinese state
3 socialism, had an effective on sensitivity and false negative reporting: rural hukou was more likely
4 to report error among those people with hypertension and have a false negative reporting. With
5 age, the rate of sensitivity went up and the rate of specificity went down. However, false reporting
6 increased with age. For diabetes, educational attainment was also an important explanatory factor,
7 and had a significant impact on sensitivity and false negative reporting, which meant the least
8 educated people were tend to underestimate their diabetes burden than their counterparts. Besides,
9 Hukou had an effective on sensitivity, specificity and correct reporting: rural hukou was more
10 likely to report error among those people with hypertension and correct report among those people
11 without diabetes. As a whole, rural hukou had a more correct reporting. Elderly people were more
12 prone to aware their diabetes, while less inclined to report their absence of disease. The false
13 reporting and false negative reporting went up significantly with age.

14 We draw three lessons from our results. First, our findings confirm the previous questions
15 that there is a big gap between self-reports and biomarker in China. China is a rapidly rising
16 developing country and is undergoing rapid population aging, which are generally associated with
17 and non-communicable and chronic conditions. Although a great progress of China in its inclusion
18 of biological and anthropometric measures of health in this and other surveys expands the
19 possibilities for biomarkers and social construction, underdiagnosis of disease is really common
20 in China. Considering the underdiagnosis of disease, the Chinese government should increase
21 awareness of disease and reassess the burden of disease. Second, self-reported data underestimate
22 the disease burden of hypertension and diabetes, and the underestimation of diabetes is greater.
23 Adding objective measurements to social survey could improve data accuracy allowing better
[Insert Running title of <72 characters]

1 understanding of socioeconomic inequalities in health. We underline the need to supplement
2 subjective health data with comprehensive and reliable biomedical measures where possible.
3 Objective measures of health, biomarkers are more valid measures of physiological function
4 “under the skin”, meaning biosocial approaches to enhance the importance of social factors in the
5 biomedical process and to intervene in social conditions that cause inequity and avoidable inequity
6 will become increasingly important (Harris and Schorpp 2018). Third, there is an urgent need to
7 provide basic health education and physical examination to citizens, to facilitate access to
8 healthcare and make focused interventions to lower the incidence and unawareness of disease in
9 China.

10 A drawback of the paper is that we used cross-sectional data, which limits causal inferences.
11 The challenge of identifying causal effects remains universal in most science research, including
12 social stratification and health research. In this regard, longitudinal data with biomarker or genetics
13 are especially useful for sorting out causal effect. For example, having baseline biomarker
14 measures prior to some social exposure or self-assessment enables researchers to identify change
15 in that biomarker response to that exposure and explore whether and to what extent age trajectories
16 of self-reported health, biomarkers and their discrepancy, depended on the educational level.
17 Another limitation is biomedical measurement is an imperfect criterion. Problems arise when
18 availability of the biomarker is differentially related to either the disease or the exposure or when
19 the specimen acquisition, storage, measurement, or ascertainment procedures differ in those with
20 the disease compared to those without the disease or outcome of interest, and the most important
21 source of confounding is the failure to identify factors that may alter the measurement of the
22 biomarker, such as metabolic factors. If biological stability is not guaranteed, the accuracy of the
23 biomarker cannot be guaranteed. Future work should also consider additional data sources and
[Insert Running title of <72 characters]

1 repeated multiple biomarker measurements to complement survey data.

2

3 **Conclusion**

4 The prevalence of hypertension and diabetes are increasing over age in China, with many old
5 people remaining undiagnosed. Self-reported hypertension and diabetes showed low sensitivity
6 (71.98% and 49.21%, respectively) but high specificity (93.71% and 98.05%, respectively). False
7 positive reporting of hypertension and diabetes were 3.85% and 1.67%, while false negative
8 reports were extremely high at 10.85% and 7.38%. Education degree, hukou, age and gender affect
9 both the specific error and the overall error of reporting hypertension and diabetes, but there are
10 some differences in the magnitude and direction. As this is the first report of undiagnosed
11 hypertension and diabetes by using four indexes and evaluate sociodemographic characteristics
12 that are correlated with misreporting in China, the results confirm self-reported conditions
13 underestimate the disease burden. Adding objective measurements into social survey could
14 improve data accuracy allowing better understanding of socioeconomic inequalities in health.

15

16

17

18

19

20

21

22

23

[Insert Running title of <72 characters]

1 **Acknowledgments**

2 We thank CHARLS team, who collected data and assisted with data access for the study.

3 **Authors Contributions**

4 DX and JW developed the study design. DH was responsible for data management, statistical
5 analysis and drafting of this paper. JW provided critical knowledge in drafting of the paper. All
6 authors read and approved final manuscript.

7 **Funding**

8 This study was supported by China Scholarship Council and U.S.-China Fulbright Program. The
9 funding body had no influence on the design and collection, analysis, and interpretation of the data
10 and the writing of the manuscript.

11 **Availability of data and materials**

12 The datasets generated and/or analyzed during the study are publicly available.
13 <http://charls.pku.edu.cn>

14 **Ethics approval**

15 The original CHARLS was approved by the Ethical Review Committee of Peking University, and
16 all participants signed informed consent at the time of participation.

17 **Consent for publication**

18 Not Applicable.

19 **Competing interests**

20 The authors declare that there is no conflict of interest regarding the publication of this paper and
21 approve the final version of the manuscript being submitted.

22

[Insert Running title of <72 characters]

1 **References**

2 Akobeng, A. K., et al. (2006). "Effect of breast feeding on risk of coeliac disease: a systematic
3 review and meta-analysis of observational studies." Archives of disease in childhood **91**(1): 39-43.

4
5 Al Shamsi, H. and A. Almutairi (2018). "Disagreement Between Self-Reporting and Objective
6 Diagnosis in Chronic Diseases Among Omanis 2008." Global Journal of Health Science **10**(5).

7
8 Alwan, A. (2011). Global status report on noncommunicable diseases 2010, World Health
9 Organization.

10
11 Baker, M., et al. (2004). "What do self-reported, objective, measures of health measure?" Journal
12 of human Resources **39**(4): 1067-1093.

13
14 Bennett, C., et al. (2007). "HbA1c as a screening tool for detection of type 2 diabetes: a systematic
15 review." Diabetic medicine **24**(4): 333-343.

16
17 Brasher, M. S., et al. (2017). "Incorporating biomarkers into the study of socio-economic status
18 and health among older adults in China." SSM-population health **3**: 577-585.

19
20 Chen, X., et al. (2019). "Venous Blood-Based Biomarkers in the China Health and Retirement
21 Longitudinal Study: Rationale, Design, and Results From the 2015 Wave." American journal of
22 epidemiology **188**(11): 1871-1877.

23
[Insert Running title of <72 characters]

1 Control, C. f. D. and Prevention (2017). "National diabetes statistics report, 2017." Atlanta, GA:
2 Centers for Disease Control and Prevention, US Department of Health and Human Services **20**.

3

4 Harris, K. M. and K. M. Schorpp (2018). "Integrating biomarkers in social stratification and health
5 research." Annual review of sociology **44**: 361-386.

6

7 Huang, G., et al. (2017). "Prevalence, awareness, treatment, and control of hypertension among
8 very elderly Chinese: results of a community-based study." Journal of the American Society of
9 Hypertension **11**(8): 503-512. e502.

10

11 Kavishe, B., et al. (2015). "High prevalence of hypertension and of risk factors for non-
12 communicable diseases (NCDs): a population based cross-sectional survey of NCDS and HIV
13 infection in Northwestern Tanzania and Southern Uganda." BMC medicine **13**(1): 126.

14

15 Kislaya, I., et al. (2019). "Differential self-report error by socioeconomic status in hypertension
16 and hypercholesterolemia: INSEF 2015 study." European journal of public health **29**(2): 273-278.

17

18 Kohler, U., et al. (2011). "Comparing coefficients of nested nonlinear probability models." The
19 Stata Journal **11**(3): 420-438.

20

21 Mayeux, R. (2004). "Biomarkers: Potential Uses and Limitations." NeuroRx **Vol**(2): 182-188..

22

[Insert Running title of <72 characters]

1 Murray, C. J. and L. C. Chen (1992). "Understanding morbidity change." The Population and
 2 Development Review: 481-503.

3

4 Newell, S. A., et al. (1999). "The accuracy of self-reported health behaviors and risk factors
 5 relating to cancer and cardiovascular disease in the general population: a critical review."
 6 American journal of preventive medicine **17**(3): 211-229.

7

8 Onur, I. and M. Velamuri (2018). "The gap between self-reported and objective measures of
 9 disease status in India." PloS one **13**(8).

10

11 Paulose-Ram, R., et al. (2017). Characteristics of US adults with hypertension who are unaware
 12 of their hypertension, 2011-2014, US Department of Health & Human Services, Centers for
 13 Disease Control and

14

15 Ritchie, H. and M. Roser (2019). "Sanitation." Our World in data.

16

17 Trevethan, R. (2017). "Sensitivity, specificity, and predictive values: foundations, pliabilitys, and
 18 pitfalls in research and practice." Frontiers in public health **5**: 307.

19

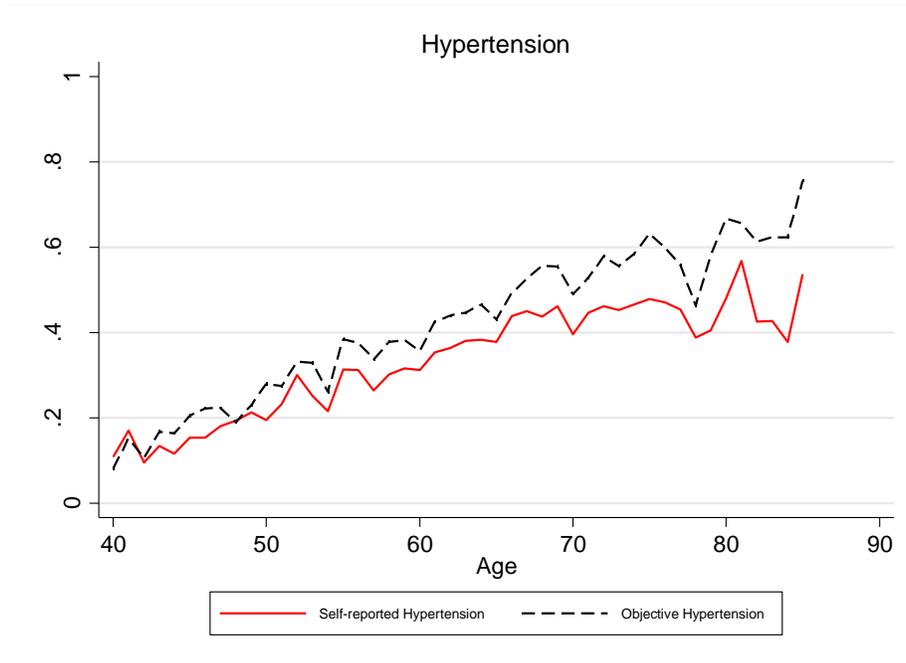
20 Zhou, M., et al. (2019). "Mortality, morbidity, and risk factors in China and its provinces, 1990–
 21 2017: a systematic analysis for the Global Burden of Disease Study 2017." The Lancet **394**(10204):
 22 1145-1158.

23

[Insert Running title of <72 characters]

1

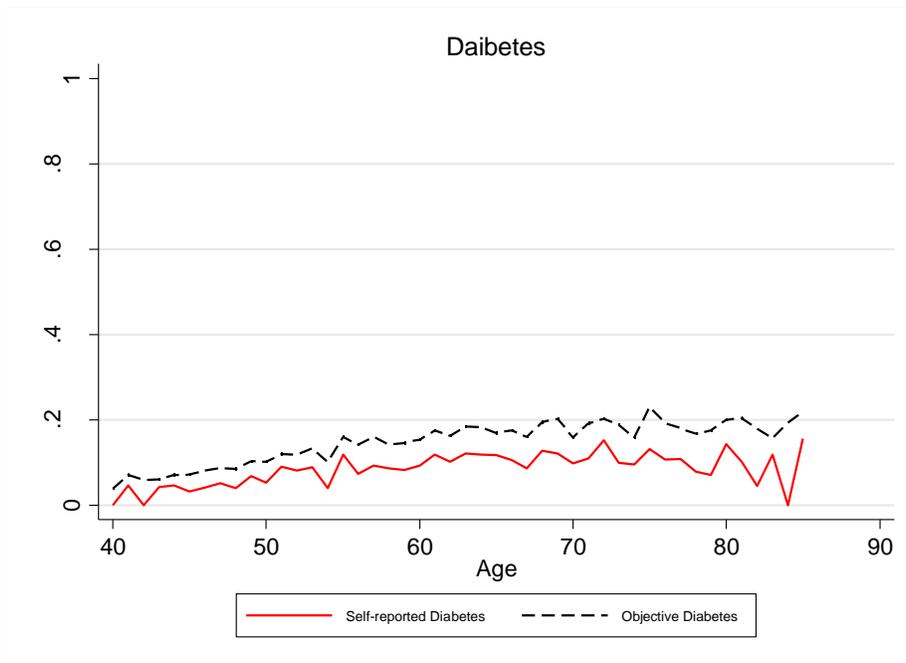
2 **Legends**



3

4 **Fig. 1** Self-reported Hypertension and Objective Hypertension by Age

5



6

7

[Insert Running title of <72 characters]

Fig. 2 Self-reported Diabetes and Objective Diabetes by Age

Tables
Table 1 Characteristics of the Sample, China Health and Retirement Longitudinal Study

Variable	Total (N=19,292)	Hypertension (N=14,457)	Diabetes (N=12,189)
Education			
Illiterate	0.244	0.254	0.256
Primary education	0.454	0.462	0.459
Secondary education and above	0.302	0.284	0.285
Hukou (Urban=0)	0.770	0.797	0.797
Drinking			
None	0.732	0.736	0.738
Less than 3 days a month	0.062	0.061	0.059
Once or 2 to 3 days a week	0.064	0.063	0.062
4 to 6 days a week or daily	0.088	0.085	0.085
Twice a day or above	0.054	0.055	0.056
Smoking	1.735	1.723	1.715
Sex (female = 0)	0.477	0.465	0.460
Age			
40-49	0.215	0.196	0.186
50-59	0.321	0.318	0.320
60-69	0.298	0.314	0.324
70-79	0.136	0.143	0.144
80 and above	0.030	0.029	0.026
Marriage (Not=0)	0.878	0.878	0.877

Table 2 Sensitivity, Specificity, False Negative Reporting and False Positive Reporting (%)

	Hypertension	Diabetes
Sensitivity	71.98	49.21
Specificity	93.71	98.05
False negative reporting	10.85	7.38
False positive reporting	3.85	1.67

[Insert Running title of <72 characters]

1
2
3
4
5

Table 3 Predicting Disagreement or False reporting of Hypertension: Odds Ratio Estimates, Logistic and Multinomial Logistic Models, with z Scores in Parentheses

	Sensitivity	Specificity	Correct Reporting		
			Correct Reporting	False Negative	False Positive
Hypertension					
Education (Illiterate=0)					
<i>Primary education</i>	1.125 (1.496)	0.872 (-1.151)	1.066 (1.019)	0.875+ (-1.900)	1.151 (1.188)
<i>Secondary and above</i>	1.211* (1.994)	0.839 (-1.252)	1.097 (1.240)	0.829* (-2.204)	1.191 (1.259)
Hukou (Urban=0)	0.732*** (-3.924)	0.945 (-0.477)	0.877* (-2.067)	1.132+ (1.696)	1.161 (1.287)
Drink (None=0)					
<i>< 3 days a /month</i>	0.988 (-0.085)	1.085 (0.424)	1.042 (0.384)	0.961 (-0.320)	0.957 (-0.229)
<i>Once or 2 to 3 days a/week</i>	0.995 (-0.040)	0.882 (-0.688)	0.894 (-1.117)	1.120 (0.975)	1.124 (0.645)
<i>4 to 6 days a/week or daily</i>	0.614*** (-4.571)	0.685* (-2.458)	0.624*** (-5.808)	1.655*** (5.523)	1.460* (2.485)
<i>Twice a day or above</i>	0.636*** (-3.569)	1.032 (0.149)	0.726** (-3.210)	1.501*** (3.717)	1.002 (0.011)
Smoke	0.920** (-3.041)	1.102* (2.329)	1.006 (0.290)	1.013 (0.516)	0.940 (-1.512)
Sex (female = 0)	0.848* (-2.100)	0.903 (-0.891)	0.840** (-2.804)	1.239** (3.010)	1.080 (0.679)
Age (40-49=0)					
<i>50-59</i>	1.345** (2.765)	0.913 (-0.712)	0.821** (-2.630)	1.379*** (3.580)	0.933 (-0.547)
<i>60-69</i>	1.870*** (5.990)	0.648*** (-3.438)	0.747*** (-3.929)	1.491*** (4.513)	1.072 (0.553)
<i>70-79</i>	1.647*** (4.280)	0.662* (-2.497)	0.628*** (-5.346)	1.957*** (6.660)	0.912 (-0.564)
<i>>=80</i>	1.242 (1.285)	0.508* (-2.170)	0.438*** (-5.995)	2.935*** (7.107)	1.013 (0.041)
Marriage (Not=0)	1.056 (0.621)	1.157 (1.015)	1.158* (2.032)	0.843* (-2.119)	0.932 (-0.501)
N	5597	8860	14457	14457	14457

6

Note: + p < 0.10; * p < 0.05; ** p < 0.01; *** p < 0.001.

[Insert Running title of <72 characters]

1

2 **Table 4** Predicting Disagreement or False reporting of Diabetes: Odds Ratio Estimates, Logistic
 3 and Multinomial Logistic Models, with z Scores in Parentheses

	Sensitivity	Specificity	Correct Reporting		
			Correct Reporting	False Negative	False Positive
Diabetes					
Education (Illiterate=0)					
<i>Primary education</i>	1.136 (1.009)	0.959 (-0.216)	1.057 (0.681)	0.930 (-0.814)	1.029 (0.149)
<i>Secondary and above</i>	1.436* (2.423)	0.745 (-1.337)	1.124 (1.177)	0.811+ (-1.916)	1.310 (1.233)
Hukou (Urban=0)	0.616*** (-4.147)	1.420* (2.050)	1.186* (2.118)	0.869 (-1.577)	0.752+ (-1.669)
Drink (None=0)					
<i>< 3 days a /month</i>	0.942 (-0.261)	0.545* (-2.311)	0.948 (-0.376)	0.892 (-0.690)	1.873* (2.389)
<i>Once or 2 to 3 days</i>	0.918 (-0.373)	0.589+ (-1.938)	0.996 (-0.028)	0.864 (-0.886)	1.732* (2.012)
<i>4 to 6 days a/week or</i>	0.516** (-3.224)	0.787 (-0.866)	0.872 (-1.139)	1.113 (0.818)	1.319 (1.000)
<i>Twice a day or above</i>	0.433** (-2.976)	0.886 (-0.346)	1.046 (0.294)	0.910 (-0.564)	1.200 (0.521)
Smoke	0.940 (-1.292)	1.122 (1.577)	1.040 (1.268)	0.976 (-0.726)	0.897 (-1.488)
Sex (female = 0)	1.005 (0.042)	1.337 (1.543)	1.117 (1.315)	0.930 (-0.780)	0.752 (-1.516)
Age (40-49=0)					
<i>50-59</i>	1.542* (2.394)	0.892 (-0.494)	0.765* (-2.505)	1.371** (2.659)	1.072 (0.304)
<i>60-69</i>	1.819*** (3.423)	0.629* (-2.074)	0.614*** (-4.707)	1.676*** (4.477)	1.456+ (1.686)
<i>70-79</i>	1.430+ (1.801)	0.521* (-2.491)	0.521*** (-5.446)	1.967*** (5.102)	1.745* (2.133)
<i>>=80</i>	1.105 (0.308)	0.440+ (-1.820)	0.456*** (-3.986)	2.228*** (3.735)	2.068 (1.617)
Marriage (Not=0)	1.277 (1.629)	0.907 (-0.431)	1.059 (0.600)	0.918 (-0.830)	1.089 (0.376)
N	1770	10419	12189	12189	12189

4

Note: + p < 0.10; * p < 0.05; ** p < 0.01; *** p < 0.001.

5

[Insert Running title of <72 characters]