

# Urvalsmetoder och Estimation 2

Sampling and Estimation 2

2011-01-24

# Seminar on web panel surveys (in Swedish)

- All of you are invited to a seminar, February 3rd, 13.00-16.30
- But the language i Swedish.  
<http://gauss.stat.su.se/wpu/>
- Send your application by e-post to:  
[joakim.malmdin@scb.se](mailto:joakim.malmdin@scb.se),
- Tell the organiser that you are a student and that Stockholm university statistical department should pay for you

# Preliminary programme

- **13.00–13.05 Välkommen**
- **13.05–13.35 Inledning** *Gösta Forsman, Trafikverket*
- **13.35–14.10 Webbpaneler i praktiken** *Henrik Kronberg, Norstat*
- **14.15–14.40 Vad säger omvärlden? ISO-standarden för accesspaneler samt skrifter från Esomar och AAPOR.** *Bengt Larsson, SMIF*
- **14.40–15.05 Kaffe**
- **15.05–15.45 Är icke-sannolikhetsurval aldrig representativa?** *Jan Wretman, Stockholms universitet*
- **15.50–16.15 Att bedöma webbpaneler och webbpanelundersökningar - numeriska mått och verbala beskrivningar.** *Meddelas senare*
- **16.15–16.30 Avslutande diskussion. Frågor och synpunkter till kommittén.** *Diskutant – meddelas senare*

# 3. Regression type estimators

## 3.1 Estimation using auxiliary variables

- A population  $U$  is given
- An auxiliary variable:  $X_i; i \in U$  is known
- Study variable  $Y_i$  unknown
- Take a sample  $S$ , observe  $Y_i; i \in S$
- Using  $Y$  from the sample and  $X$  from the population try to find a good estimator

# Some methods to use auxiliary information

- We will now concentrate on the estimation phase
- The sample will be assumed taken and we will for simplicity assume SRS  
(Everything works for other sampling schemes too, but more complicated).
- Design-based approach
- Estimator of the population total

# Some estimation techniques using auxiliary variables

- Difference estimators
- Ratio estimators
- Regression estimators
- Generalised regression estimators
- Prediction estimators

## 3.2 Difference estimators

- Suppose that we can make a prior guess of the unknown  $Y_i$ -value for all units using the auxiliary variables. Here we call the guess  $X_i$ .
- For example last years value or last years value plus inflation.
- Look at the differences:  $E_i = Y_i - X_i$
- Estimate the total difference  $T_e$  by  $t_e = \sum_S E_i$  as in SRS
- Estimate the total  $T_y$  by  $t_{yD} = T_x + t_e$ , where  $T_x$  is known
- Estimate variance accordingly  $\text{Var}^*(t_{yD}) = \text{Var}^*(t_e) = ((N-n)/N) \sum_S (E_i - t_e)^2 / (n-1)$

# Are difference estimators good?

- The variance is
$$\text{Var}(t_{yD}) = \text{Var}(T_x + t_e) = \text{Var}(t_e) = \text{Var}(t_y - t_x) = \text{Var}(t_y) + \text{Var}(t_x) - 2\text{Cov}(t_y, t_x)$$

- We gain if:  $2\text{Cov}(t_y, t_x) > \text{Var}(t_x)$
- If the guess is good we have  $\text{Var}(t_y) \sim \text{Var}(t_x)$   
then we gain if  $\rho(t_y, t_x) = \rho(Y, X) > 1/2$

(The reverse martingale property allows us to omit the correction for finite population)

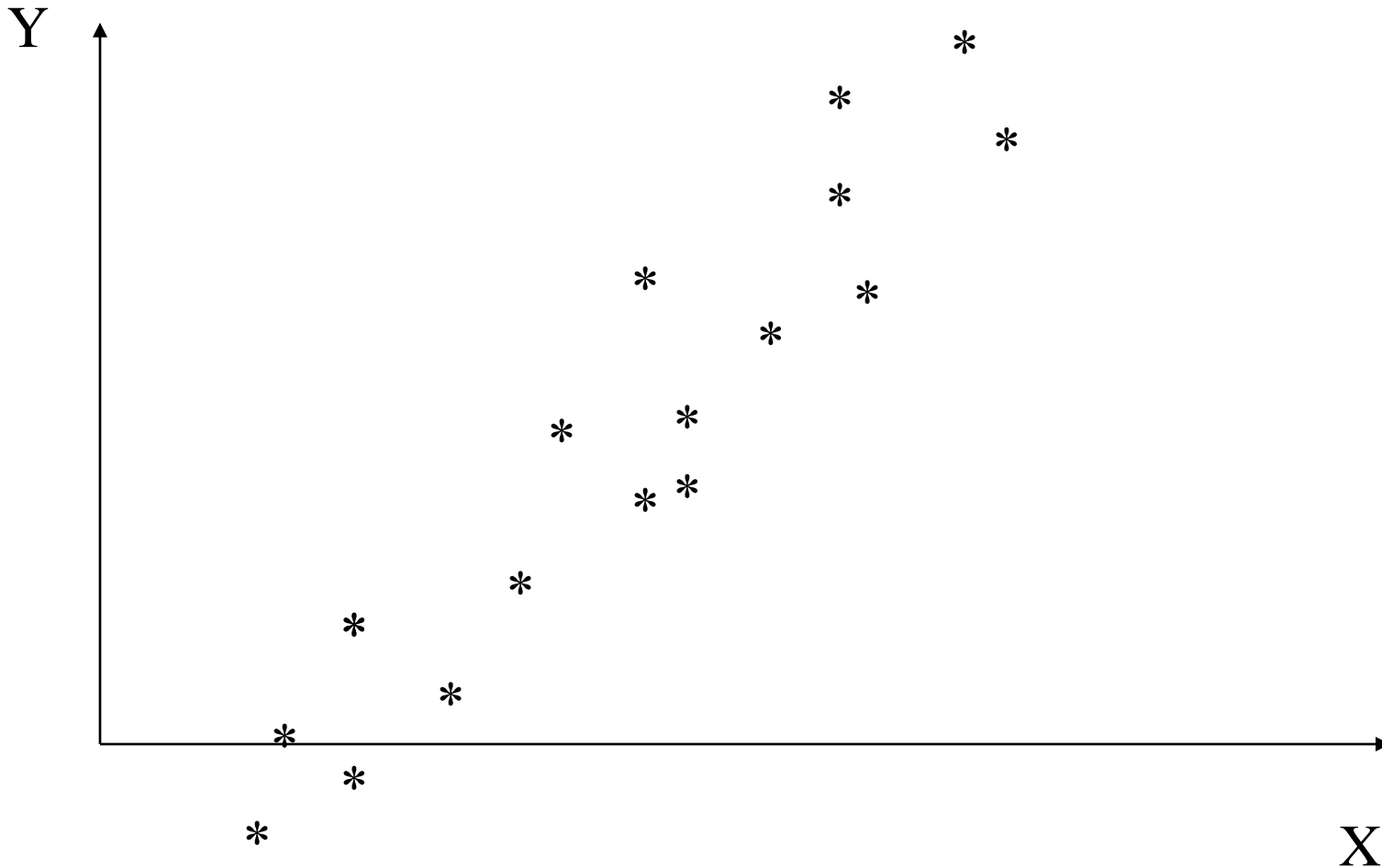
- Otherwise if  $\rho(Y, X) > \sigma_x / 2\sigma_y$



# Are difference estimators good?

- Not used so often as a basic approach
- Often used in secondary analyses and longitudinal approaches
- One problem is that it often needs small changes i.e.  $X$  has the same level as  $Y$ .
- Another problem is that there often are large variances in  $X$
- Difference estimators are good as a building blocks for other estimators

The recommended situation is when the data lies on a straight line with slope one (or slightly more than one)



# Optimal

- $\text{Var}(t_y - at_x) = \text{Var}(t_y) + a^2\text{Var}(t_x) - 2a\text{Cov}(t_y, t_x) = \text{Var}(t_y) + a^2\text{Var}(t_x) - 2a\rho(Y, X)\text{Var}(t_y)\text{Var}(t_x)$
- This is minimised when  $a = \rho(Y, X)$  (if the variances are equal), but we do not know the correlation or if the variances are equal.
- Thus the recommendation is use only for high covariances and similar variances.
- Otherwise (see regression estimators below with estimated parameters).

## 3.3 Ratio estimators

- In difference estimators we looked at the difference  $T_Y = T_X + T_{Y-X}$  and estimated it by  $t_{yD} = T_X + t_{y-x} = T_X + t_y - t_x$
- Here we look at the ratio instead  
$$t_{yR} = T_X * t_y / t_x$$
- Sensible only for positive variables

# Are ratio estimators good?

- Works best if both the mean and the variance of  $Y$  given  $X$  increase linearly with  $X$
- May be slightly biased.  
(Problem since  $t_y/t_x$  is a convex function in  $t_x$  )
- To compute the approximate bias and variance we need the theorem Gauss approximation in two dimensions.

But first

## Gauss approximation in one dimension

(or Taylor linearisation or the  $\Delta$ -method)

- If  $X$  is a random variable with variance and  $f$  a twice differentiable function

$$\text{Var}(f(X)) \approx \text{Var}(f(m_X) + (X - m_X)f'(m_X)) =$$

$$= (f'(m_X))^2 \text{Var}(X)$$

$$E(f(X)) \approx$$

$$\approx E(f(m_X) + (X - m_X)f'(m_X) + (X - m_X)^2 f''(m_X)/2) =$$

$$= f(m_X) + f''(m_X) \text{var}(X)/2$$

Gauss approximation in two dimensions: If  $X$  and  $Y$  are random variables with variances and  $f$  is twice continuously differentiable then

$$Var(f(X, Y)) \approx$$

$$\left(\frac{\partial f(m_x, m_y)}{\partial x}\right)^2 Var(X) + \left(\frac{\partial f(m_x, m_y)}{\partial y}\right)^2 Var(Y) +$$

$$2\left(\frac{\partial f(m_x, m_y)}{\partial x}\right)\left(\frac{\partial f(m_x, m_y)}{\partial y}\right)Cov(X, Y)$$

and

$$E(f(X, Y)) \approx f(m_x, m_y) + \left(\frac{\partial^2 f(m_x, m_y)}{\partial x^2}\right)Var(X) / 2 +$$

$$\left(\frac{\partial^2 f(m_x, m_y)}{\partial y^2}\right)Var(Y) / 2 + \left(\frac{\partial^2 f(m_x, m_y)}{\partial x \partial y}\right)Cov(X, Y)$$

We use the second expression on

$$t_{yR} = T_X * t_y / t_x \text{ and get } \text{Bias}(t_{yR}) =$$

$$\begin{aligned} &\approx T_x \left( \frac{\text{Var}(t_x) E(t_y)}{E^3(t_x)} + \frac{\text{Cov}(t_y, t_x)}{E^2(t_x)} \right) \\ &\approx \frac{N}{n} \left( \frac{\sigma_x^2 m_y}{m_x^2} - 2 \frac{\sigma_{xy}}{m_x} \right) \frac{N-n}{N} \end{aligned}$$

The bias is thus of order  $N/n$  and may be large in particular if the variance coefficient of  $X$  ( $\text{Var}^{1/2}(X)/E(X)$ ) is large (i.e.  $E(X)$  close to 0).



- We use the first expression on  $t_{yR} = T_X * t_y / t_x$  and get  $\text{Var}(t_{yR}) =$

$$\begin{aligned} &\approx T_x^2 \left( \frac{\text{Var}(t_y)}{E^2(t_x)} + \frac{\text{Var}(t_x)E^2(t_y)}{E^4(t_x)} - 2 \frac{\text{Cov}(t_y, t_x)E(t_y)}{E^3(t_x)} \right) \\ &\approx \frac{N^2 m_y^2}{n} \left( \frac{\sigma_y^2}{m_y^2} + \frac{\sigma_x^2}{m_x^2} - 2 \frac{\sigma_{xy}}{m_x m_y} \right) \left( \frac{N-n}{N} \right) \end{aligned}$$

Equivalently and easier to remember:

$$\text{RelVar}(t_{yR}) \sim \text{RelVar}(t_y) + \text{RelVar}(t_x) - 2\text{RelCov}(t_y, t_x)$$

Where RelVar stands for the relative variance or coefficient of variation and RelCov for relative covariance

(The Var for difference estimator is replaced by RelVar)

# Variance estimator

- One may use the above expression for the variance and replace all unknown parameters by their estimates

$m_x$  by  $t_x/n$ ;       $m_y$  by  $t_y/n$ ;       $\sigma_y^2$  by  $\Sigma_s (y - t_y/n)^2 / (n-1)$ ;  
 $\sigma_x^2$  by  $\Sigma_s (x - t_x/n)^2 / (n-1)$ ;  
 and  $\sigma_{xy}^2$  by  $\Sigma_s (y - t_y/n)(x - t_x/n) / (n-1)$ :

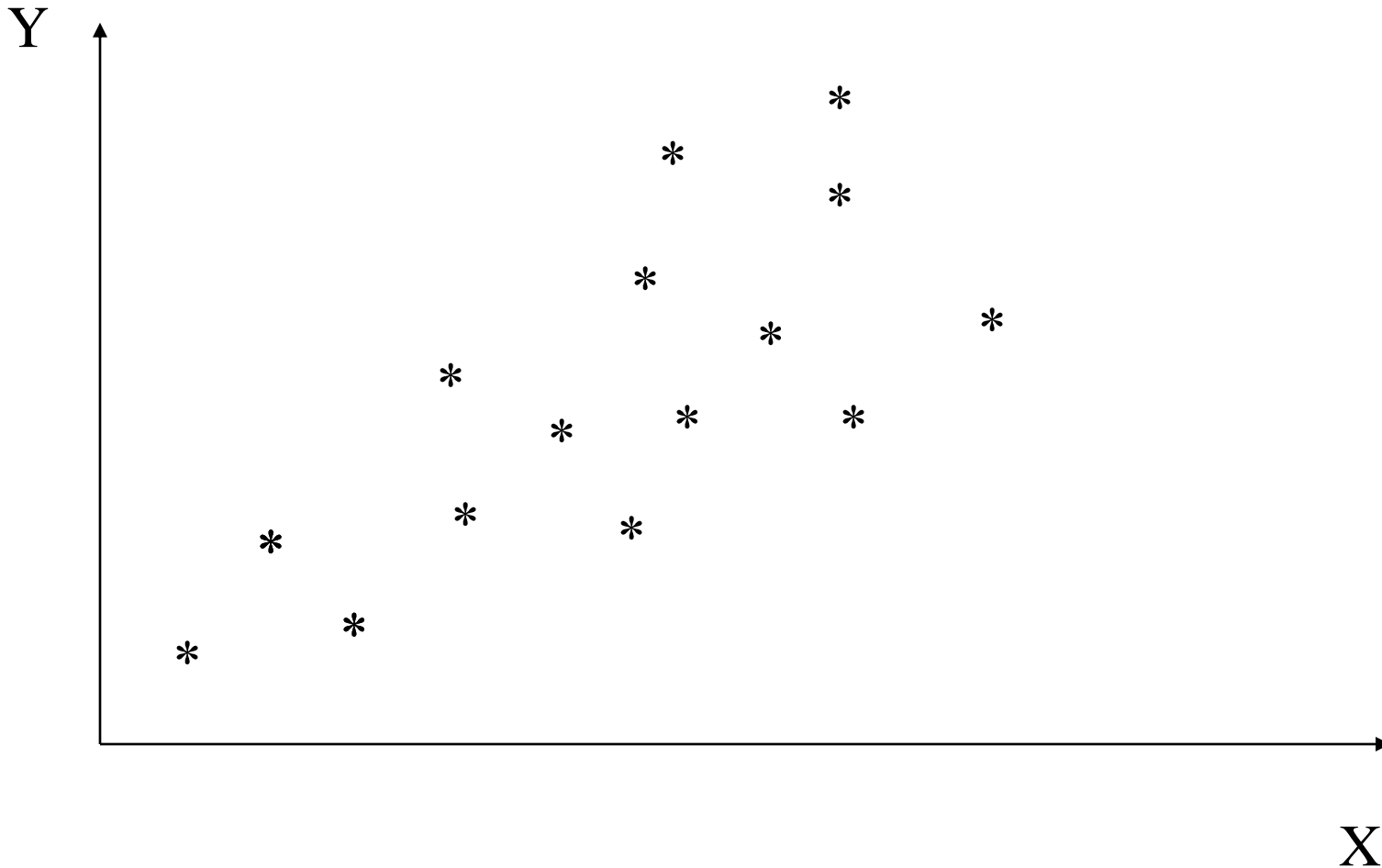
- We will show another expression later  

$$\text{Var}^*(t_{Y,R}) = N(N-n)/n \Sigma_s (E_i^* - \Sigma_s E_i^* / n)^2 / (n-1)$$
 with  $E_i^* = Y_i - (t_y / t_x) X_i$ .

# Are ratio estimators good?

- Much more often used (than difference estimator).
- One gains if the correlation is larger than  $\frac{1}{2}$  (if the relative variances (= variation coefficients) are the same)
- Multiplicative relations are more often encountered in practice than additive (All size dependent variables)
- Variance often increases with size, often linearly.

The best situation is when the data lies on a straight line through the origin and has a linearly increasing variance (or curves slightly upwards)



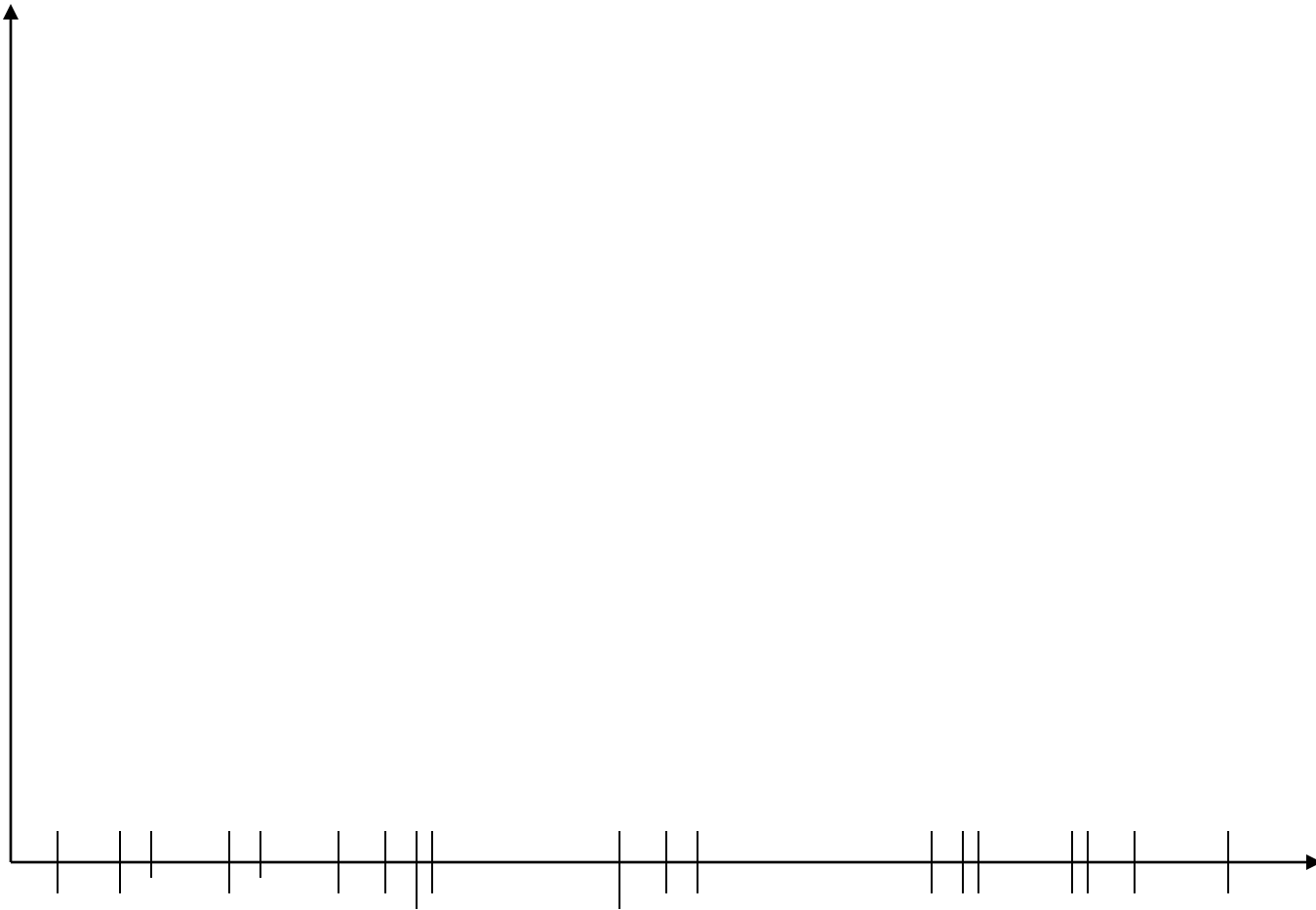
## 3.4 Regression estimator

- Same set up: A population  $U$  and a known auxiliary variable:  $X_i; i \in U$
- Take a sample  $S$ , observe a study variable  $Y_i; i \in S$
- Idea: Using the sample, find a relation between  $X$  and  $Y$ . Try to use this relation and that  $X$  is known in the estimation phase

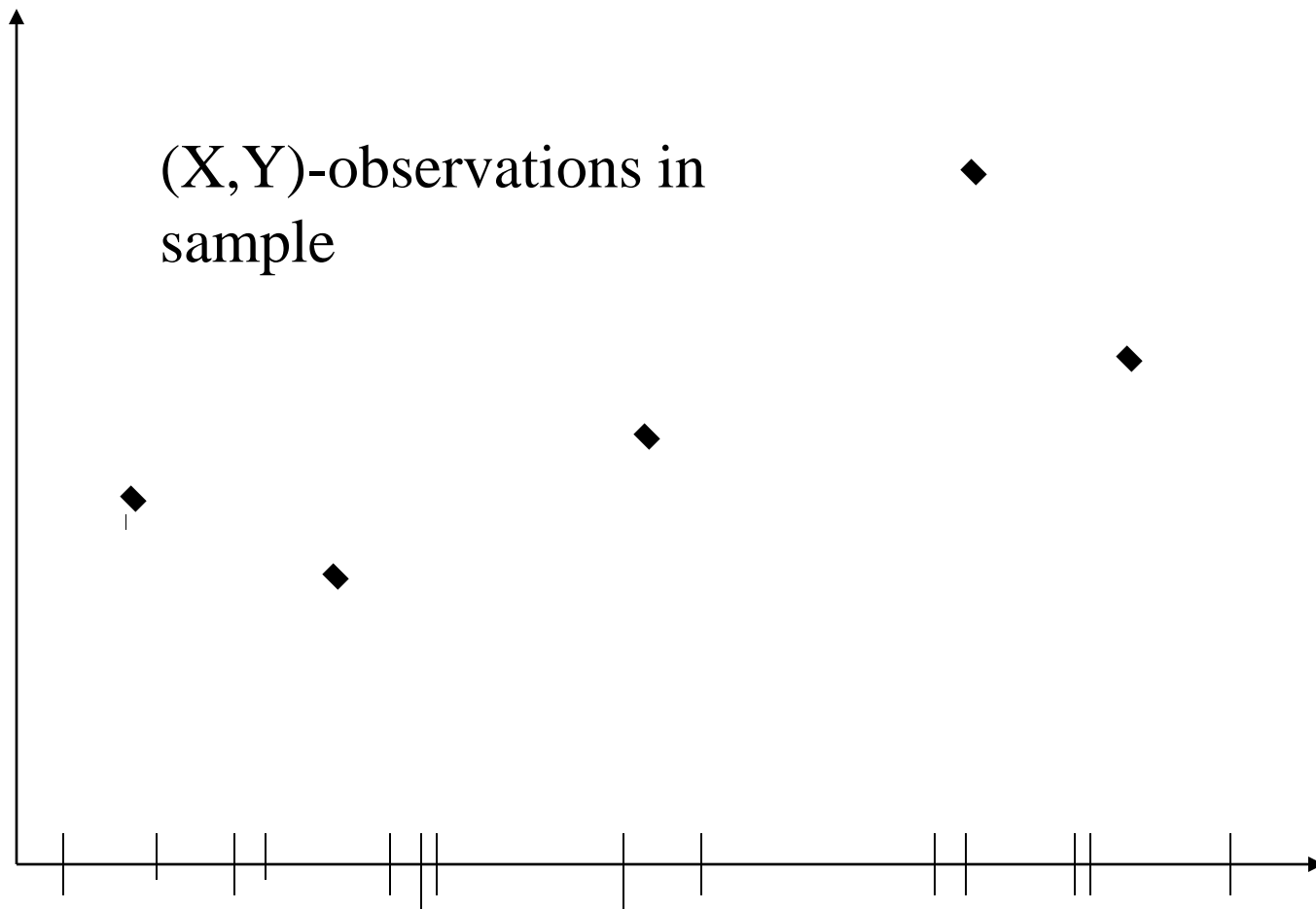
# Regression estimator

- Simple linear relation: Estimate  $a$  and  $b$  in the relation  $Y_i = a + bX_i + \varepsilon_i$  by  $a^*$  and  $b^*$
- Predict  $Y_i^* = a^* + b^* X_i$ ;  $i \in U-S$
- Estimate population total by
$$T_Y^* = \sum_S Y_i + \sum_{U-S} Y_i^*$$
$$= a^* N + b^* T_X$$

The second equality holds if  $a$  and  $b$  are estimated by Simple Linear Regression

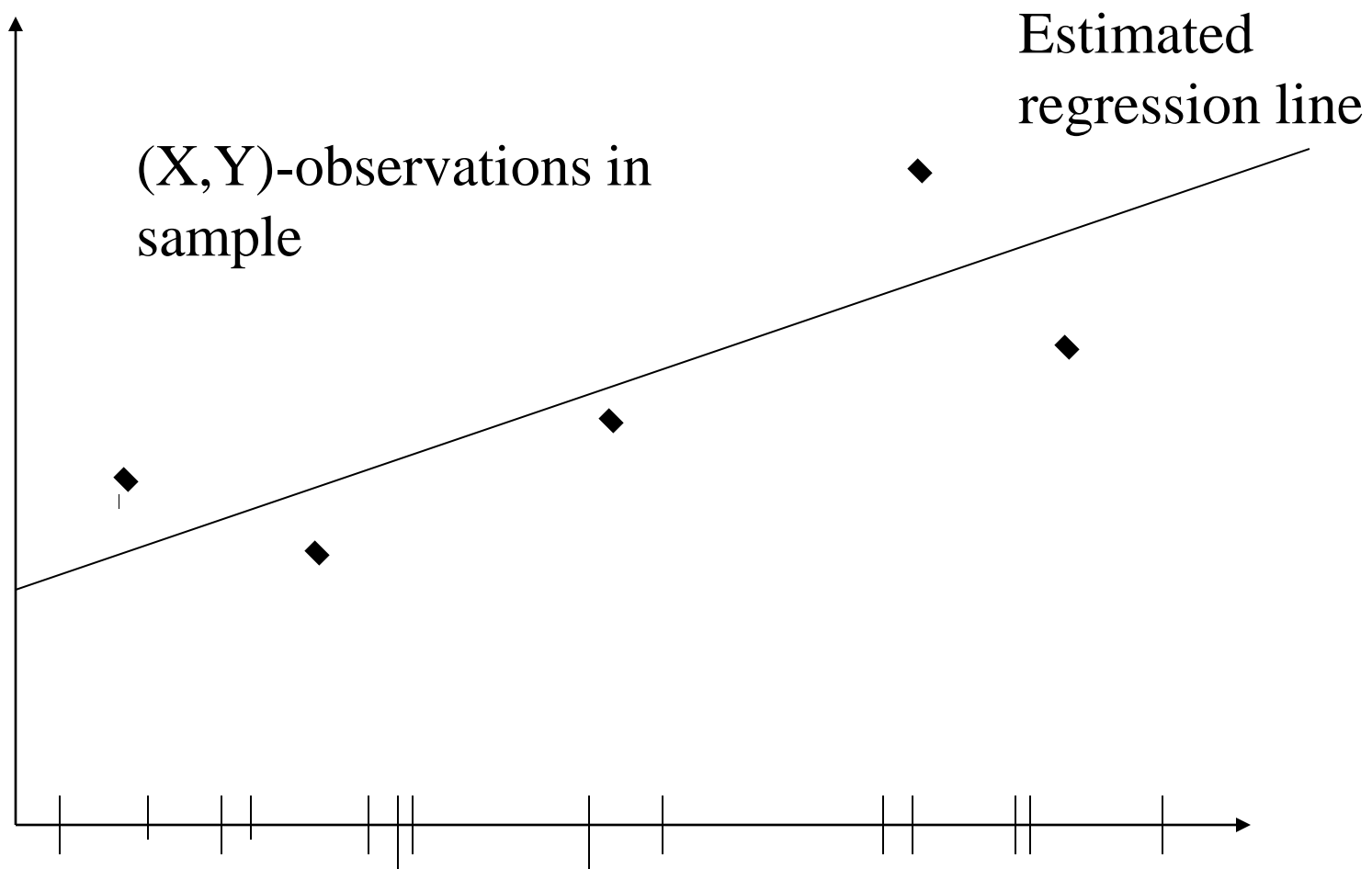


X-values in frame

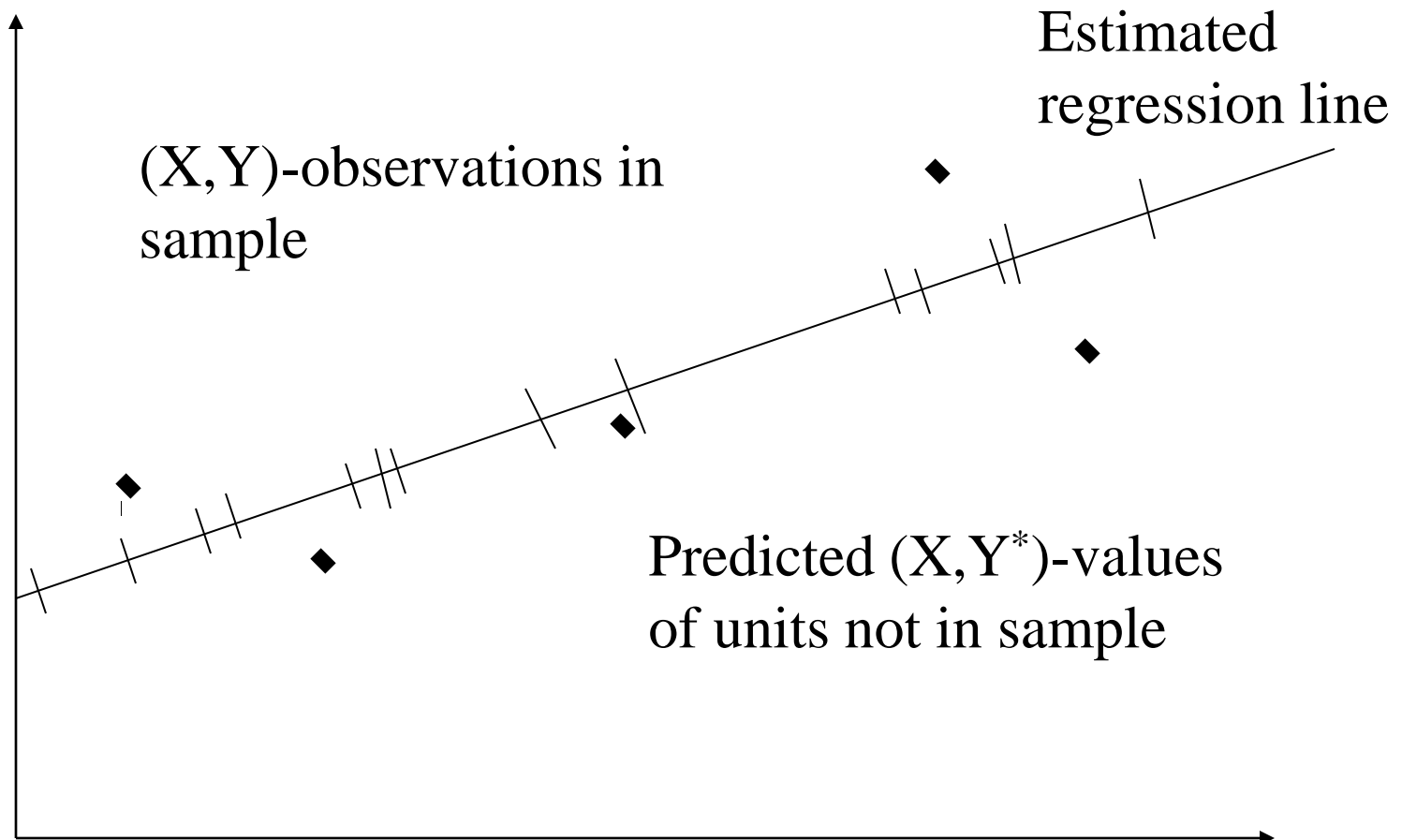


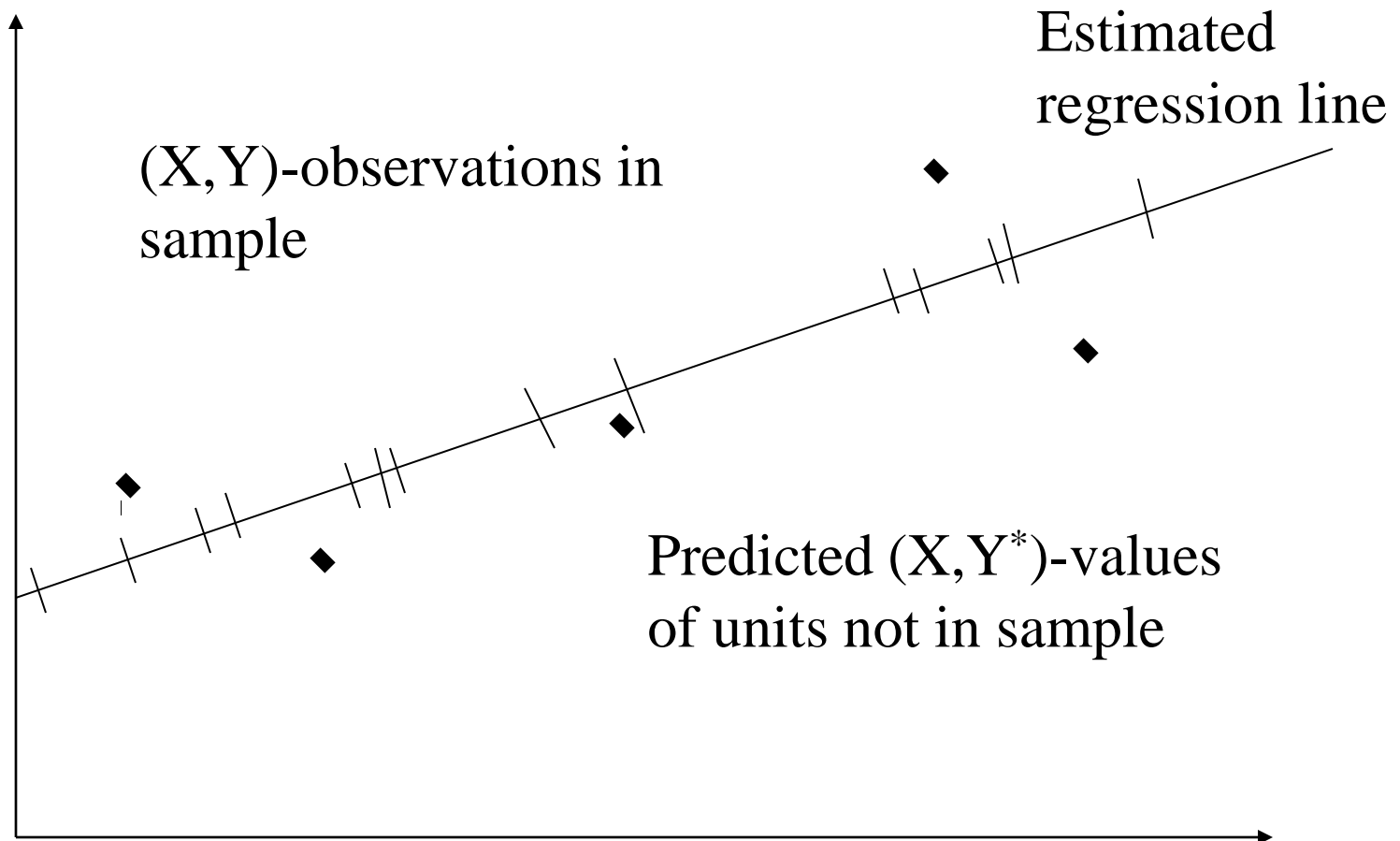
X-values in frame (except sample)





X-values in frame (except sample)





Predict the total by summing all  $Y$  and  $Y^*$ -values in population

# Regression estimator

## Further comments

- Is it sensible with Simple Linear Regression? Parameters can be estimated in other ways, and the method works then too.
- One method is the Asymptotically Optimal Regression Estimator. This approach minimises the asymptotic error variance (Montanari, (1987), ISR, 55). Based on the estimated covariance matrix of  $(t_y, t_x)$ .
$$t_{y, \text{Mont}} = t_y + (\text{Cov}^*(t_y, t_x) / \text{Var}^*(t_x)) (T_x - t_x)$$
- But usually only marginally better and may be much worse for small sample sizes (same for SRS)

# Variance and variance estimation

- Suppose first we know  $a$  and  $b$
- Write  $E_i = Y_i - (a + b X_i)$
- Then  $T_Y^* = a N + b T_X + (N/n) \sum_S E_i$
- In that case the variance of the estimator is just the variance of  $t_e = (N/n) \sum_S E_i$ , the standard estimate from SRS with  $E$  instead of  $Y$ . (cf difference estimators)
- $\text{Var}^*(t_e) = N(N-n)/n \sum_U (E_i - \sum_U E_i/N)^2/(N-1)$
- And the variance estimator is  
 $\text{Var}^*(t_e) = N(N-n)/n \sum_S (E_i - \sum_S E_i/n)^2/(n-1)$

# Variance and variance estimator (cont.)

- The variance with unknown  $a$  and  $b$  is approximately the same  

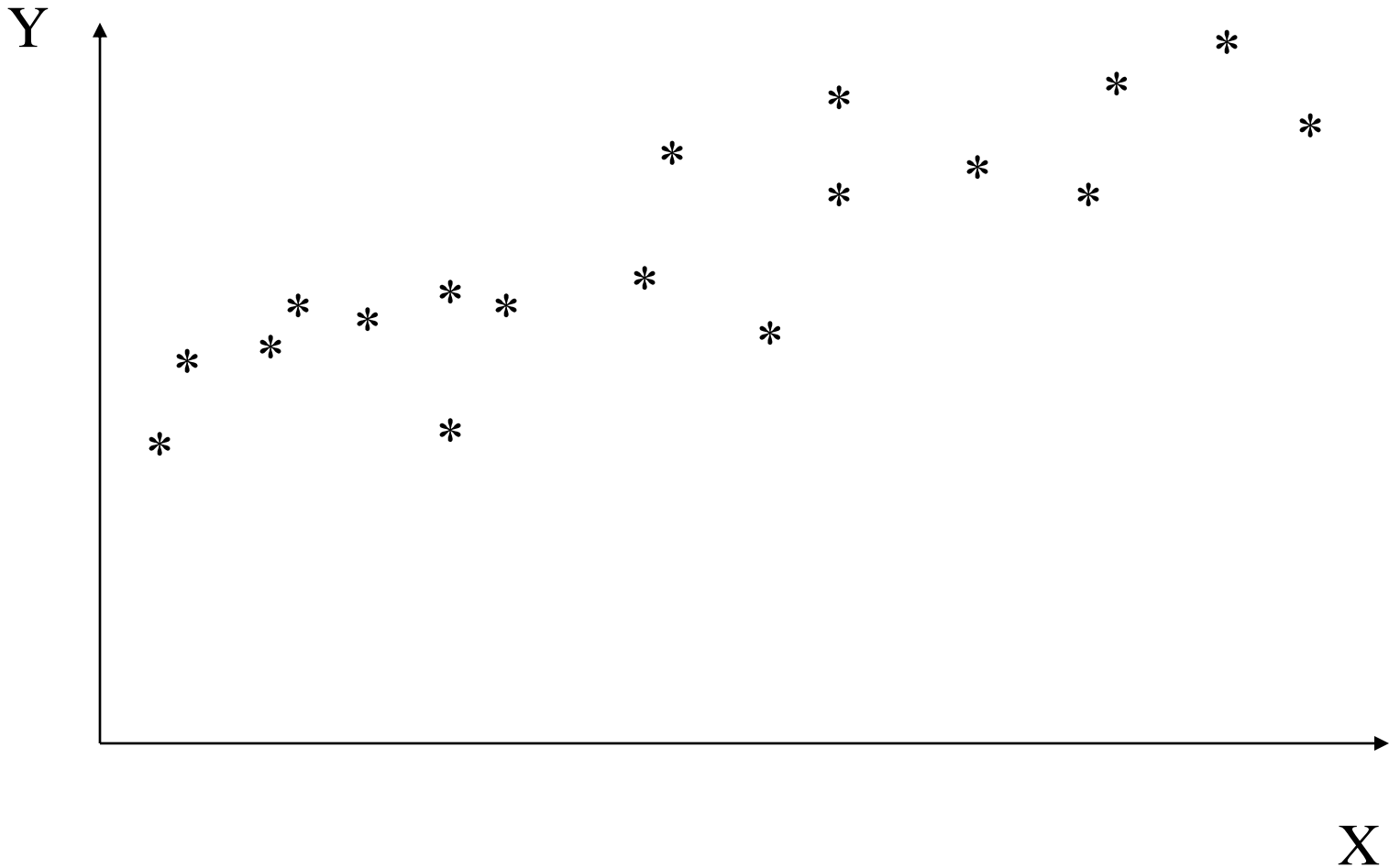
$$\text{Var}(t_{Y,\text{reg}}) \sim N(N-n)/n \sum_U (E_i - \sum_U E_i / N)^2 / (N-1)$$
- In variance estimation one usually just replaces  $a$  and  $b$  by their estimates ( $a^*$  and  $b^*$ ).
- The variance estimator (with known  $a$  and  $b$ ) was  

$$\text{Var}^*(t_e) = N(N-n)/n \sum_S (E_i - \sum_S E_i / n)^2 / (n-1)$$
- The variance can be estimated by replacing  $E_i = Y_i - (a + b X_i)$  by  $E_i^* = Y_i - (a^* + b^* X_i)$  i.e.
- $$\text{Var}^*(t_{Y,\text{reg}}) = N(N-n)/n \sum_S (E_i^* - \sum_S E_i^* / n)^2 / (n-1)$$

# Are regression estimates good?

- Improves the asymptotic variance as soon as the correlation is different from zero (compared to mean estimate, difference or ratio estimator. The decrease depends on the the "explained variance"  $R^2$ )
- May be worse for small sample sizes
- The procedure works even if the "true" relation is not linear. The estimate may become slightly biased  
(Since the estimate  $b^*$  and the mean  $t_x$  may be dependent or since the regression estimator can't be computed, if all sampled x-values have the same value. Extremely unlikely but enough to destroy exact unbiasedness)
- Care must be taken when the sample is taken with varying inclusion probabilities. (Not today's topic). (The best line may be different in sample and the remaining population)

The best situation is when the data lies on a straight line with constant variance





## 3.5 General regression estimates (GREG) 1

- Previously: only one (continuous) X-variable
- What is said can easily be generalised to several auxiliary variables using multiple linear regression
- X-values can be discrete, categorical (dummy-variables) or derived (e.g.:  $x^2$  or  $x_1 * x_2$ )
- Everything holds with  $E_i = Y_i - (a + \sum_j b_j * X_{ij})$  in the approximate variance estimation expression

# Weighted general regression estimates

- Different weighting in the regression (optimally weights should be proportional to  $\text{Var}(y_i|x_i)$ . "Model-assisted approach")

Ordinary:  $b^* =$

$$\Sigma(y_i - \bar{y})(x_i - \bar{x}) / \Sigma(x_i - \bar{x})^2$$

Weighted:  $b^* =$

$$\Sigma((y_i - \bar{y})(x_i - \bar{x}) / \text{Var}(y_i|x_i)) / \Sigma((x_i - \bar{x})^2 \text{Var}(y_i|x_i))$$

if intercept is unknown

- Other weights are sometimes used
- Weights may cause problems when the relation is not linear. The slope is estimated for points where the weight is high, but is used for all points.
- A general regression estimator is always unbiased but this does not hold for a weighted GREG-estimator.

- Difference estimators can be thought of as regression estimators with known slope (equal to one)
- Ratio estimators can be thought of as regression estimators, where the intercept is known (0) and the weights are proportional to  $x_i$   
(Holds for Poisson distribution and often for economic variables from e.g. firms)

# Are general regression estimators good?

- Asymptotically never worse than simple estimators or regression estimators with a subset of the auxiliary variables
- With many auxiliary variables the random error increases due to estimation problems of the regression coefficients. (In particular if many auxiliary variables are irrelevant or are only weakly related to  $Y$ ).
- The variance estimator underestimates the true variance
- Avoid regression estimator when there are "outliers" in the population with a set of  $X$ -values which differ much from the others or if there are influential outliers in the sample.

## An alternative expression for the variance estimator of the ratio estimator

- As we said the ratio estimator can be viewed as a regression estimator with known intercept  $a=0$ .
- The variance estimator of the regr. estimator was  
$$\text{Var}^*(t_{Y,\text{reg}}) = N(N-n)/n \sum_S (E_i^* - \sum_S E_i^*/n)^2/(n-1)$$
with  $E_i^* = Y_i - (a^* + b^* X_i)$
- Thus a variance estimator for the ratio estimator is  
$$\text{Var}^*(t_{Y,\text{reg}}) = N(N-n)/n \sum_S (E_i^* - \sum_S E_i^*/n)^2/(n-1)$$
with  $E_i^* = Y_i - (0 + b^* X_i)$  where  
$$b^* = t_y / t_x = \sum_S Y_i / \sum_S X_i$$

(This is the alternative expression we mentioned above)

## 3.6 Prediction estimates

- For regression estimates we wrote

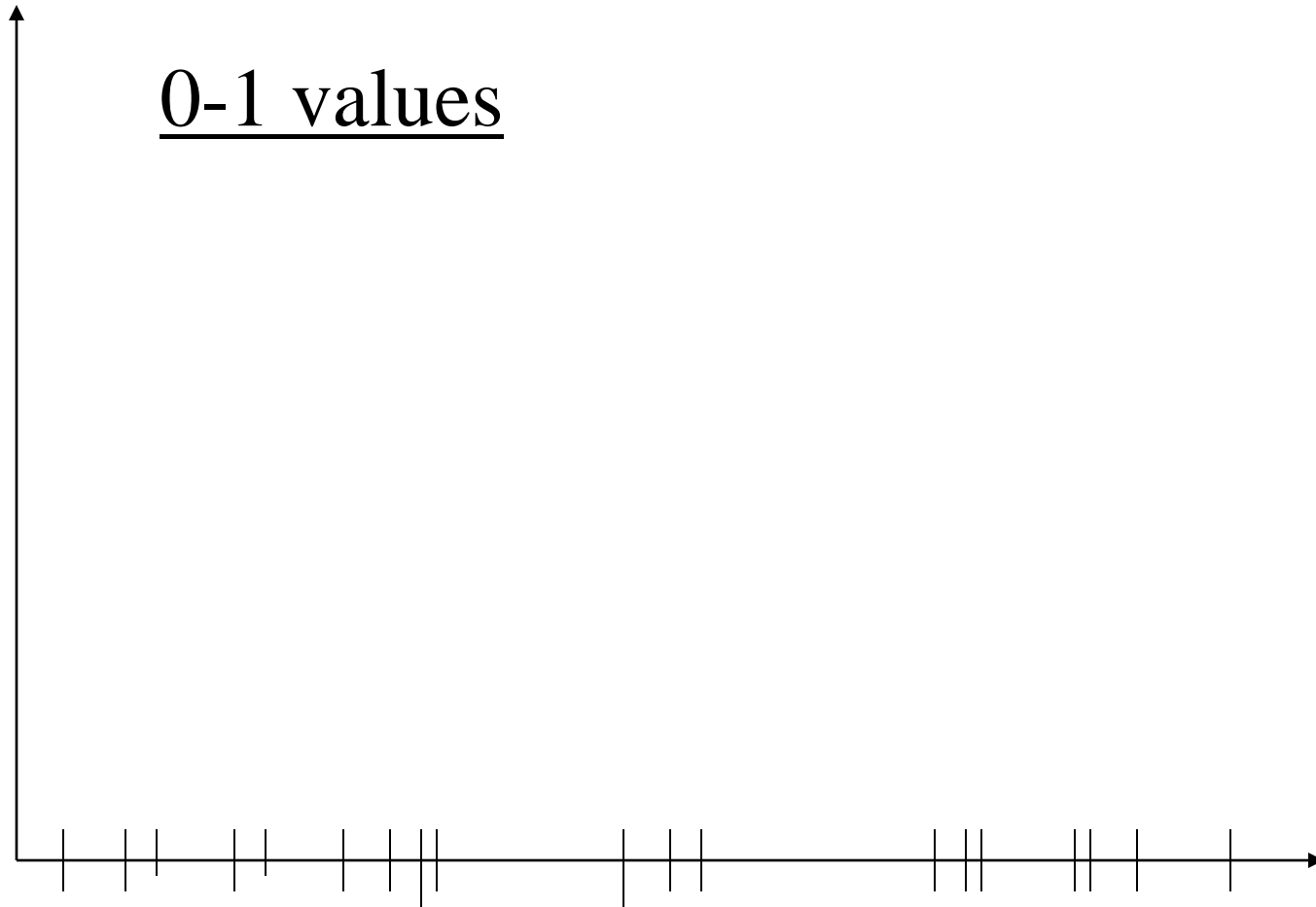
$$T_Y^* = \sum_S Y_i + \sum_{U-S} Y_i^*$$

where  $Y_i^*$  was the ordinary best linear predictor (BLUE)

- One may try to use the same formula with other predictors.
- Sensible if the predictor is good. But a bias that is unimportant for predicting one single unit may become disastrous when summed over the whole population.

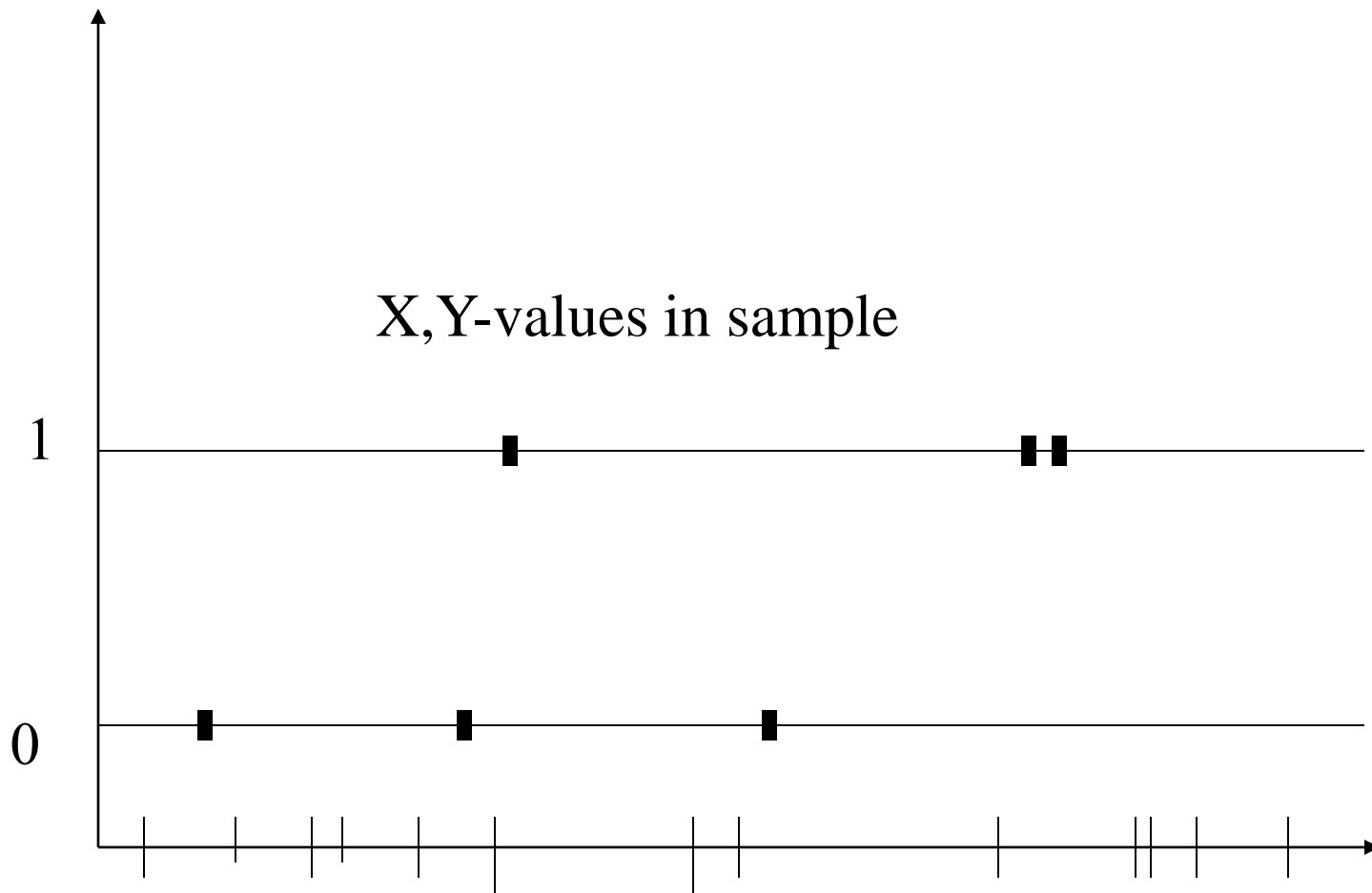
# Prediction estimates

- Idea of a difference estimator: "Suppose that we have prior guesses of the values for all units,  $X_i$ . These guesses are called  $Y_i^*$ "
- Replace the guess  $X_i$  by  $Y_i^*$ , in the full population
- Look at the difference  $E_i = Y_i - Y_i^*$
- Estimate the total difference  $T_e$  by  $t_e$
- Estimate the total  $T_y$  by  $t_{yPred} = T_{Y^*} + t_e$
- Estimate variance accordingly  $\text{Var}^*(t_{yPred}) = \text{Var}^*(t_e)$   
(asymptotically valid under mild restrictions)
- Example: 0-1-variables, proportions

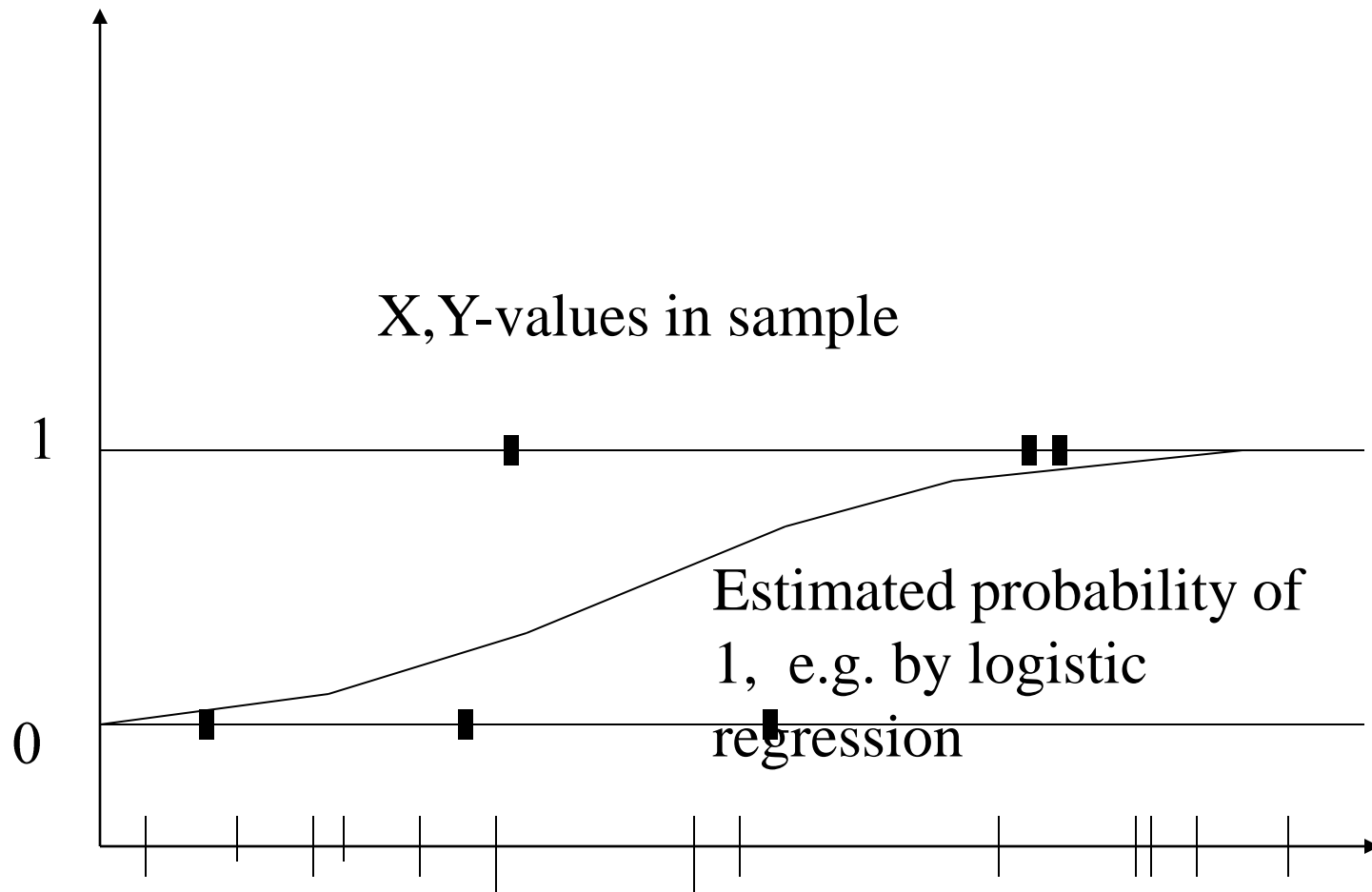


X-values in frame

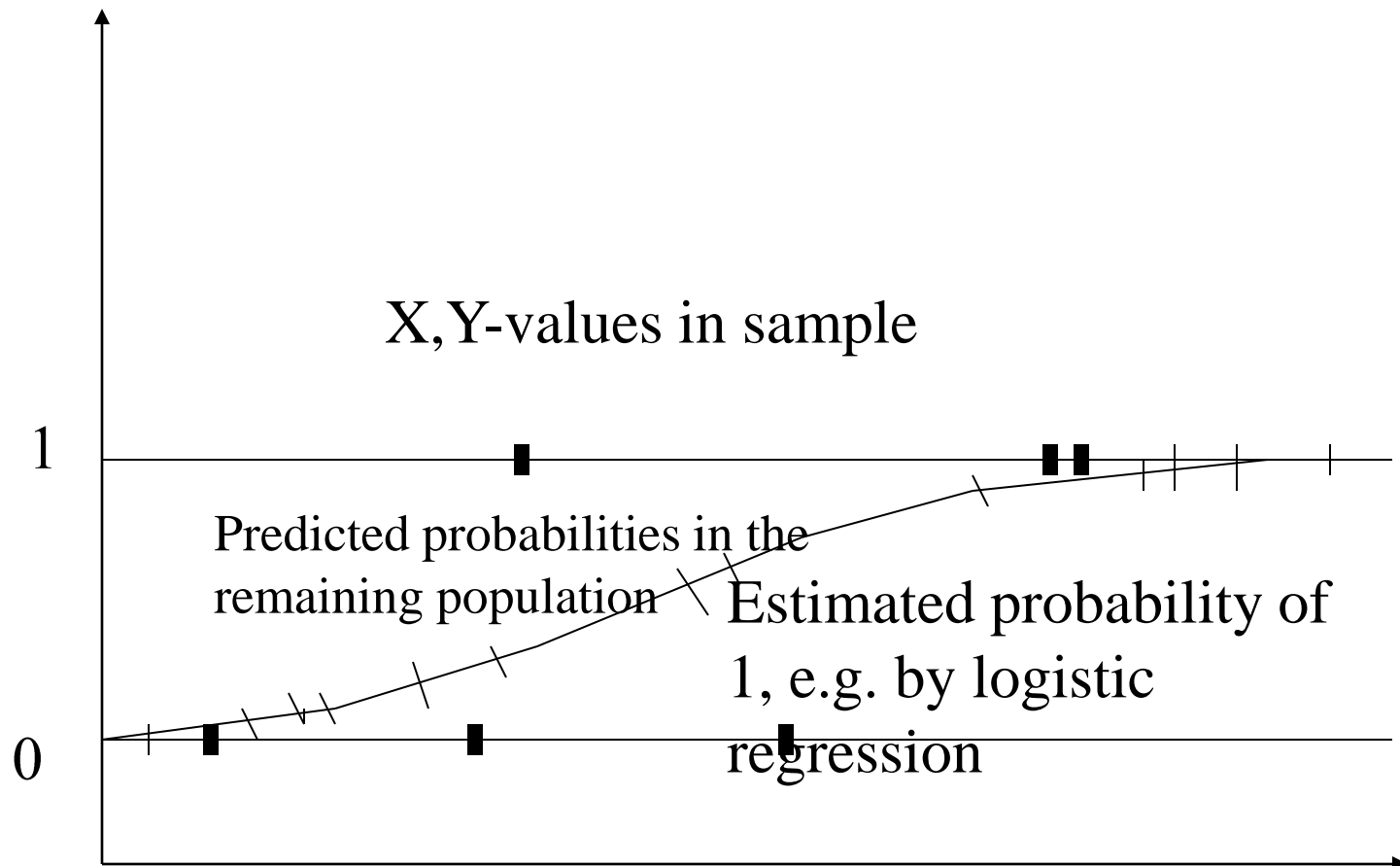




X-values in frame except in sample



X-values in frame except in sample



Estimate sum of predicted and true Y-values

- The prediction is not necessarily unbiased from a design-based perspective (in particular if the chosen (logistic) model does not hold)
- They ought thus to be corrected, to be almost unbiased for any population and large samples  

$$t_{Y,\text{pred}} = \sum_U Y_i^* + (N/n) \sum_S (Y_i - Y_i^*)$$
- Note that  $Y_i^*$  is computed also for units in the sample.
- If one believed in the logistic model and used modelbased inference the second term would not be needed.

- Prediction estimators can also be seen as approximate difference estimators
- Treat the prediction  $Y_i^*$  as the auxiliary variable (as we did for the regression estimator) and write

$$t_{Y,\text{pred}} = \sum_U Y_i^* + (N/n) \sum_S E_i$$

$$\text{where } E_i^* = Y_i - Y_i^*$$

- The variance of this can be estimated by
- $$\text{Var}^*(t_{Y,\text{pred}}) = N(N-n)/n \sum_S (E_i^* - \sum_S E_i^*/n)^2/(n-1)$$
- (exactly as for the regression and ratio estimators)

# Are prediction estimators good?

Similar comments as for (general) regression estimators:

- Asymptotically never worse than ordinary mean
- But if the sample is small and the estimated function uncertain one may lose efficiency.
- The better the model fit, the better the estimator
- Be careful with varying inclusion probabilities, outliers in the X-population and influential observations in the sample
- ...