

Lecture 25:

Regression

## Prediction

Suppose  $Y$  is some r.v.

What's the "best" guess about  $Y$ 's value?

By "best" we mean minimizes mean squared error.

## Prediction

Suppose  $Y$  is some r.v.

What's the "best" guess about  $Y$ 's value?

By "best" we mean minimizes mean squared error.

Theorem: the minimizer of  $\underline{E[(Y - \alpha)^2]}$  is  
 $\alpha = EY.$

the minimum mean squared error is



## Prediction

Suppose  $Y$  is some r.v.

What's the "best" guess about  $Y$ 's value?

By "best" we mean minimizes mean squared error.

Theorem: the minimizer of  $E[(Y - \alpha)^2]$  is  
 $\alpha = EY$ .

the minimum mean squared error is



Proof:  $E[(Y - \alpha)^2] = E[Y^2] - 2\alpha E[Y] + \alpha^2$

This is a parabola in  $\alpha$   $a\alpha^2 + b\alpha + c$

with minimal value at  $\alpha = \frac{-b}{2a} = \frac{2EY}{2} = EY$ .

## Prediction

Suppose  $X, Y$  are some r.v.

Suppose we learned  $X$ 's value.

What's the best guess for  $Y$ 's value?

Recap

# Multiple Random Variables

Joint Distribution: If  $X$  and  $Y$  are two r.v.s over the same probability space then their joint distribution is defined as

$$\{(a, b, \Pr[X=a, Y=b]) : a \in \text{range}(X), b \in \text{range}(Y)\}$$

## Marginal Distributions

Marginal for  $X$ :  $\Pr[X=a] = \sum_{b \in \text{range}(Y)} \Pr[X=a, Y=b]$

Marginal for  $Y$ :  $\Pr[Y=b] = \sum_{a \in \text{range}(X)} \Pr[X=a, Y=b]$

Recap

# Multiple Random Variables

Joint Distribution:

$$\{(a, b, \Pr[X=a, Y=b]) : a \in \text{range}(X), b \in \text{range}(Y)\}$$

Marginal for X:  $\Pr[X=a] = \sum_{b \in \text{range}(Y)} \Pr[X=a, Y=b]$

Marginal for Y:  $\Pr[Y=b] = \sum_{a \in \text{range}(X)} \Pr[X=a, Y=b]$

Example:

		Y		
	X	1	2	3
1		0	0.1	0.2
2		0.3	0	0
3		0.1	0.2	0.1

$$\Pr[X=1] =$$

$$\Pr[Y=3] =$$

$$\Pr[X=1 | Y=3] =$$

## Conditional Expectation

Def'n: Let  $X$  and  $Y$  be r.v.s over  $\Omega$ .

The conditional expectation of  $Y$  given  $X$  is defined as

$$E[Y|X] = g(X)$$

where

$$g(x) = E[Y|X=x] = \sum_y y \cdot \Pr[Y=y | X=x].$$



## Conditional Expectation

Def'n: Let  $X$  and  $Y$  be r.v.s over  $\Omega$ .

The conditional expectation of  $Y$  given  $X$  is defined as

$$E[Y|X] = g(X)$$

where

$$g(x) = E[Y|X=x] = \sum_y y \cdot \Pr[Y=y | X=x].$$

Note:  $E[Y|X]$  is a r.v. that is a fuc. of  $X$ .

$\forall x$ :  $E[Y|X=x]$  is a number.

## Properties of Conditional Expectation

$$E[Y|X=x] = \sum_y y \cdot \Pr[Y=y|X=x]$$

1. If  $X, Y$  indep.  $\Rightarrow E[Y|X] = E[Y]$ .

2.  $E[aY+b|X] = a \cdot E[Y|X] + b$

3.  $\forall$  fnc  $h(\cdot)$   $E[h(X)Y|X] = h(X) \cdot E[Y|X]$ .

4.  $E[E[Y|X]] = E[Y]$

5.  $\forall$  fnc  $h(\cdot)$   $E[h(X)E[Y|X]] = E[h(X) \cdot Y]$ .

# Properties of Conditional Expectation

$$E[Y|X=x] = \sum_y y \cdot \Pr[Y=y|X=x]$$

1. If  $X, Y$  indep.  $\Rightarrow E[Y|X] = E[Y]$ .

2.  $E[aY+b|X] = a \cdot E[Y|X] + b$

3.  $\forall$  fnc  $h(\cdot)$   $E[h(X)Y|X] = h(X) \cdot E[Y|X]$ .

4.  $E[E[Y|X]] = E[Y]$

5.  $\forall$  fnc  $h(\cdot)$   $E[h(X)E[Y|X]] = E[h(X) \cdot Y]$

---

Proof:

(1) true, by definition. (2) true, by linearity of expectation.

(3) For any  $x$ ,  $E[h(X) \cdot Y|X=x] = h(x) \cdot E[Y|X=x]$ .

$$\begin{aligned} (4) \quad E_x[E_Y[Y|X]] &= \sum_x \Pr[X=x] \cdot \sum_y y \cdot \Pr[Y=y|X=x] \\ &= \sum_y y \cdot \underbrace{\sum_x \Pr[X=x] \cdot \Pr[Y=y|X=x]}_{\Pr[Y=y]} = E[Y]. \end{aligned}$$

5. For any  $h(\cdot)$

$$\mathbb{E}_x [h(x) \mathbb{E}_Y [Y|x]] = \mathbb{E}_{X,Y} [h(x) Y].$$

5. For any  $h(\cdot)$   $\mathbb{E}_X [h(x) \mathbb{E}_Y [Y|x]] = \mathbb{E}_{X,Y} [h(x)Y]$ .

Proof:

$$\mathbb{E}_X [h(x) \cdot \mathbb{E}_Y [Y|x]] =$$

$$= \sum_x P_r [X=x] \cdot h(x) \cdot \sum_y P_r [Y=y|X=x] \cdot y.$$

$$= \sum_{x,y} h(x) \cdot y \cdot P_r [X=x, Y=y]$$

$$= \mathbb{E}_{X,Y} [h(x) \cdot Y].$$

Corollary: For all  $h(\cdot)$   $E[(Y - E[Y|X]) \cdot h(X)] = 0$

Corollary: For all  $h(\cdot)$   $E[(Y - E[Y|X]) \cdot h(X)] = 0$

Proof:

$$E[(Y - E[Y|X])h(X)] =$$

$$= E[Y \cdot h(X)] - E[E[Y|X]h(X)]$$

$$= E[Y \cdot h(X)] - E[Y \cdot h(X)].$$

Theorem: Let  $X, Y$  be two r.v.s over  $\Omega$ .

The best predictor of  $Y$  from  $X$  (minimizes mean squared error)

is  $g(x) = E[Y|X]$ .



Theorem: Let  $X, Y$  be two r.v.s over  $\Omega$ .

The best predictor of  $Y$  from  $X$  (minimizes mean squared error)

is  $g(x) = E[Y|X]$ .

Proof: Let  $h(X)$  be any function of  $X$ .

$$\begin{aligned} E[(Y-h(x))^2] &= E[(Y-g(x) + g(x)-h(x))^2] \\ &= E[(Y-g(x))^2] + 2E[(Y-g(x)) \cdot (g(x)-h(x))] + E[(g(x)-h(x))^2] \\ &\geq E[(Y-g(x))^2] \end{aligned}$$

Theorem: Let  $X, Y$  be two r.v.s over  $\Omega$ .

The best predictor of  $Y$  from  $X$  (minimizes mean squared error)

is  $g(x) = E[Y|X]$ .

Proof: Let  $h(X)$  be any function of  $X$ .

$$\begin{aligned} E[(Y-h(X))^2] &= E[(Y-g(x) + g(x)-h(x))^2] \\ &= E[(Y-g(x))^2] + 2 \underbrace{E[(Y-g(x)) \cdot (g(x)-h(x))]}_0 + \underbrace{E[(g(x)-h(x))^2]}_{\geq 0} \\ &\geq E[(Y-g(x))^2] \end{aligned}$$

Theorem: Let  $X, Y$  be two r.v.s over  $\Omega$ .

The best predictor of  $Y$  from  $X$  (minimizes mean squared error) is  $g(x) = E[Y|X]$ .

Proof: Let  $h(X)$  be any function of  $X$ .

$$\begin{aligned} E[(Y-h(X))^2] &= E[(Y-g(x) + g(x)-h(x))^2] \\ &= E[(Y-g(x))^2] + 2 \underbrace{E[(Y-g(x))(g(x)-h(x))]}_0 + \underbrace{E[(g(x)-h(x))^2]}_{\geq 0} \\ &\geq E[(Y-g(x))^2] \end{aligned}$$

Alternative Proof: For any  $x \in \text{range}(X)$ , given  $X=x$

the best predictor of  $Y$  is  $E[Y|X=x]$  from earlier.

# Linear Regression

So far, we've seen:

- If we want to guess  $Y$  without knowing anything else

best guess is  $EY$ .

- If we make some observation  $X$  related to  $Y$ :

best guess is  $g(x) = E[Y|X]$ .

# Linear Regression

So far, we've seen:

- If we want to guess  $Y$  without knowing anything else  
best guess is  $EY$ .

- If we make some observation  $X$  related to  $Y$ :

best guess is  $g(x) = E[Y|X]$ .

The latter is optimal but can be complicated.

What if we want a simpler func of  $X$  explaining  $Y$ .

For example: a linear function.

## Motivation: Statistics

In real-life applications, we don't necessarily know the joint dist. of  $X, Y$ .

We can get estimates for  $EX, EY$ , etc. from observations.

## Motivation: Statistics

In real-life applications, we don't necessarily know the joint dist. of  $X, Y$ .

We can get estimates for  $EX$ ,  $EY$ , etc. from observations.

To estimate  $E[Y|X=x]$  we need

samples such that  $X=x$ , but typically we'll have few or no such examples.

## Motivation: Statistics

In real-life applications, we don't necessarily know the joint dist. of  $X, Y$ .

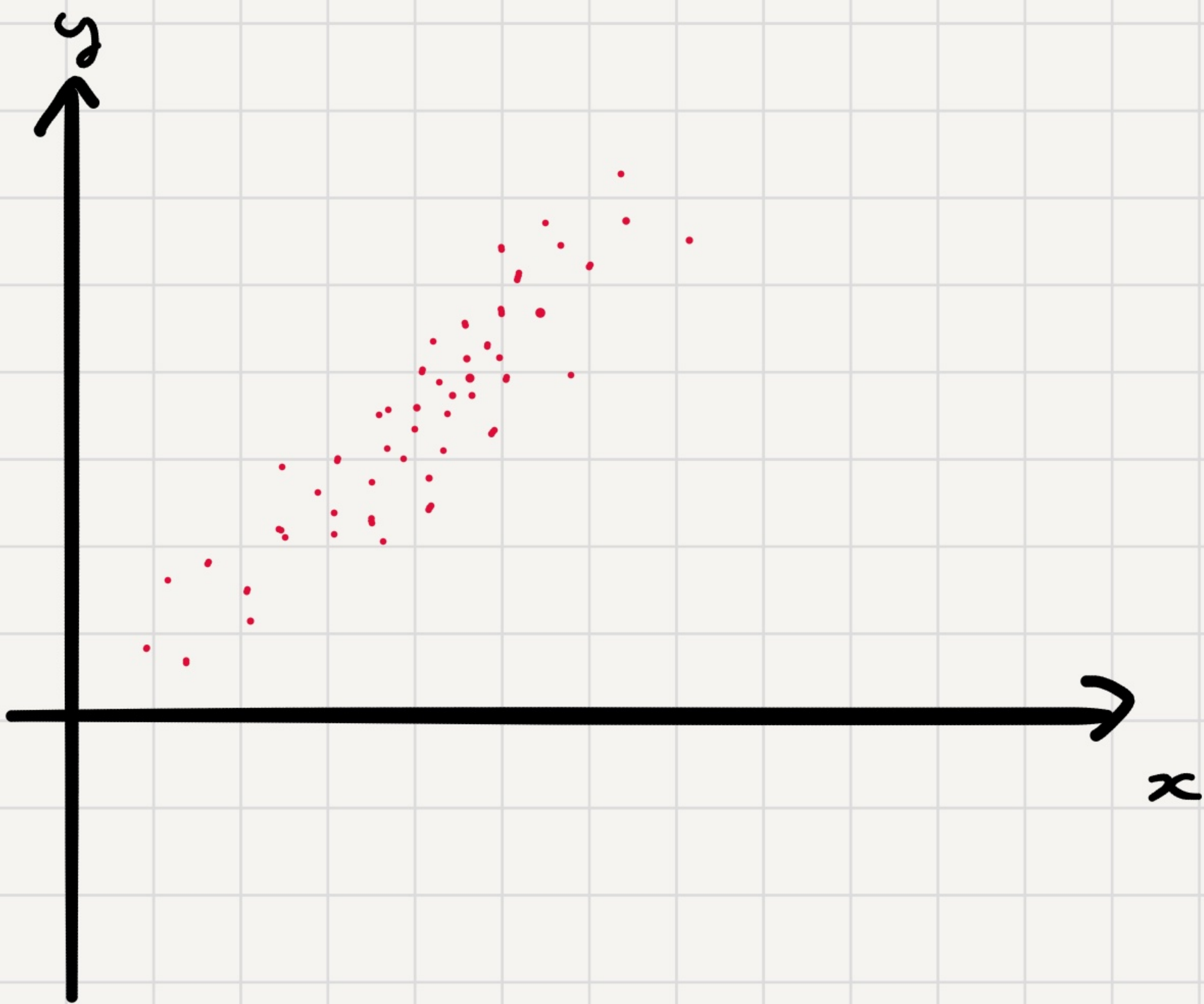
We can get estimates for  $\mathbb{E}X$ ,  $\mathbb{E}Y$ , etc. from observations.

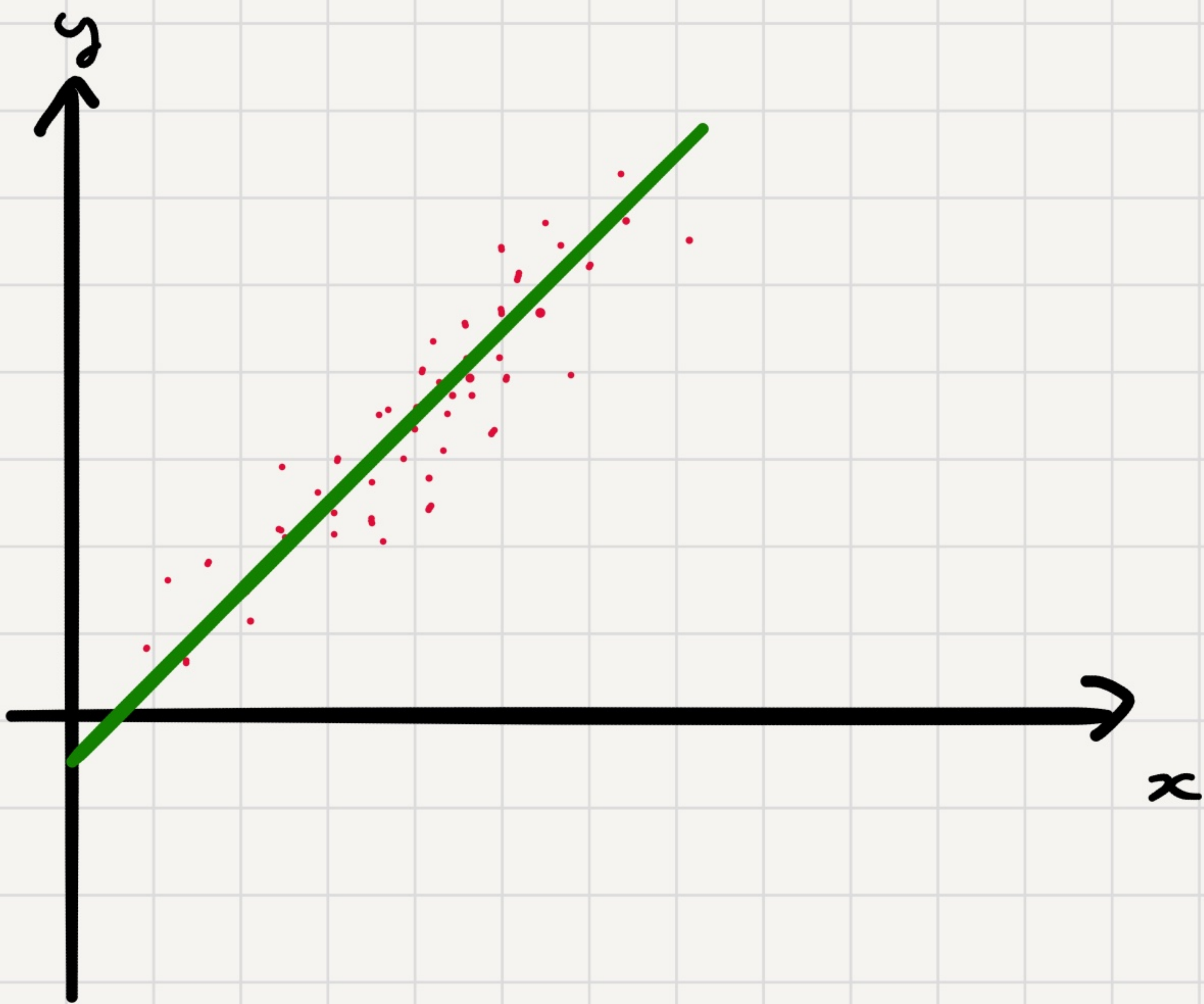
To estimate  $\pi_{1,1}$

samples such as  $(1,1)$ , but typically we'll have few or no such examples.

For a simpler model,  $\alpha X + \beta$ , we can use all the samples to get a good estimate of the two parameters:  $\alpha, \beta$ .







Theorem: The best linear predictor of  $Y$  as a function of  $X$  is  $E[Y] + \frac{\text{cov}(X, Y)}{\text{var}(X)} \cdot (X - E[X])$ .

Theorem: The best linear predictor of  $Y$  as a function of  $X$  is  $E[Y] + \frac{\text{cov}(X, Y)}{\text{var}(X)} \cdot (X - E[X])$ .

Proof: We'll first consider predicting  $\hat{Y} = Y - E[Y]$   
from  $\hat{X} = X - E[X]$ .

Theorem: The best linear predictor of  $Y$  as a function of  $X$  is  $E[Y] + \frac{\text{cov}(X, Y)}{\text{var}(X)} \cdot (X - E[X])$ .

Proof: We'll first consider predicting  $\hat{Y} = Y - E[Y]$  from  $\hat{X} = X - E[X]$ .

$$\begin{aligned} \min_{\alpha, \beta} E[(\hat{Y} - (\alpha \hat{X} + \beta))^2] &= \\ &= \min_{\alpha, \beta} [E[\hat{Y}^2] - 2E[\hat{Y}(\alpha \hat{X} + \beta)] + E[(\alpha \hat{X} + \beta)^2]] \end{aligned}$$

Theorem: The best linear predictor of  $Y$  as a function of  $X$  is  $E[Y] + \frac{\text{cov}(X, Y)}{\text{var}(X)} \cdot (X - E[X])$ .

Proof: We'll first consider predicting  $\hat{Y} = Y - E[Y]$  from  $\hat{X} = X - E[X]$ .

$$\begin{aligned} \min_{\alpha, \beta} E[(\hat{Y} - (\alpha \hat{X} + \beta))^2] &= \\ &= \min_{\alpha, \beta} [E[\hat{Y}^2] - 2E[\hat{Y}(\alpha \hat{X} + \beta)] + E[(\alpha \hat{X} + \beta)^2]] \\ &= \min_{\alpha, \beta} [E[\hat{Y}^2] - 2\alpha E[\hat{X}\hat{Y}] + \alpha^2 E[\hat{X}^2] + \beta^2] \end{aligned}$$

Theorem: The best linear predictor of  $Y$  as a function of  $X$  is  $E[Y] + \frac{\text{cov}(X, Y)}{\text{var}(X)} \cdot (X - E[X])$ .

Proof: We'll first consider predicting  $\hat{Y} = Y - E[Y]$  from  $\hat{X} = X - E[X]$ .

$$\begin{aligned} \min_{\alpha, \beta} E[(\hat{Y} - (\alpha \hat{X} + \beta))^2] &= \\ &= \min_{\alpha, \beta} [E[\hat{Y}^2] - 2E[\hat{Y}(\alpha \hat{X} + \beta)] + E[(\alpha \hat{X} + \beta)^2]] \\ &= \min_{\alpha, \beta} [E[\hat{Y}^2] - 2\alpha E[\hat{X}\hat{Y}] + \alpha^2 E[\hat{X}^2] + \beta^2] \end{aligned}$$

Solution:  $\beta = 0$ ,  $\alpha = \frac{E[\hat{X}\hat{Y}]}{E[\hat{X}^2]} = \frac{\text{cov}(X, Y)}{\text{var}(X)}$

Theorem: The best linear predictor of  $Y$  as a function of  $X$  is  $E[Y] + \frac{\text{cov}(X, Y)}{\text{var}(X)} \cdot (X - E[X])$ .

Proof: We'll first consider predicting  $\hat{Y} = Y - E[Y]$  from  $\hat{X} = X - E[X]$ .

$$\begin{aligned} \min_{\alpha, \beta} E[(\hat{Y} - (\alpha \hat{X} + \beta))^2] &= \\ &= \min_{\alpha, \beta} [E[\hat{Y}^2] - 2E[\hat{Y}(\alpha \hat{X} + \beta)] + E[(\alpha \hat{X} + \beta)^2]] \\ &= \min_{\alpha, \beta} [E[\hat{Y}^2] - 2\alpha E[\hat{X}\hat{Y}] + \alpha^2 E[\hat{X}^2] + \beta^2] \end{aligned}$$

Solution:  $\beta = 0$ ,  $\alpha = \frac{E[\hat{X}\hat{Y}]}{E[\hat{X}^2]} = \frac{\text{cov}(X, Y)}{\text{var}(X)}$

$\Rightarrow$  Best linear predictor for  $Y - E[Y]$  is  $\alpha \cdot (X - E[X])$   
 $\Rightarrow$  " " " " "  $Y$  is  $E[Y] + \alpha \cdot (X - E[X])$



Theorem: The best linear predictor of  $Y$  as a function of  $X$  is  $l(x) = E[Y] + \frac{\text{cov}(x, Y)}{\text{var}(x)} \cdot (x - E[x])$ .

Corollary: The minimum  $E[(Y - l(x))^2] = \text{Var}(Y) \cdot (1 - \text{Corr}(X, Y)^2)$ .

Theorem: The best linear predictor of  $Y$  as a function of  $X$  is  $EY + (x - EX) \cdot \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$ .

Corollary: The minimum  $E[(Y - l(x))^2] = \text{Var}(Y) \cdot (1 - \text{Corr}(X, Y)^2)$ .

Proof: By the theorem  $l(x) = EY + (x - EX) \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$

$$E[(Y - l(x))^2] = E\left[\left((Y - EY) - (x - EX) \frac{\text{Cov}(X, Y)}{\text{Var}(X)}\right)^2\right]$$

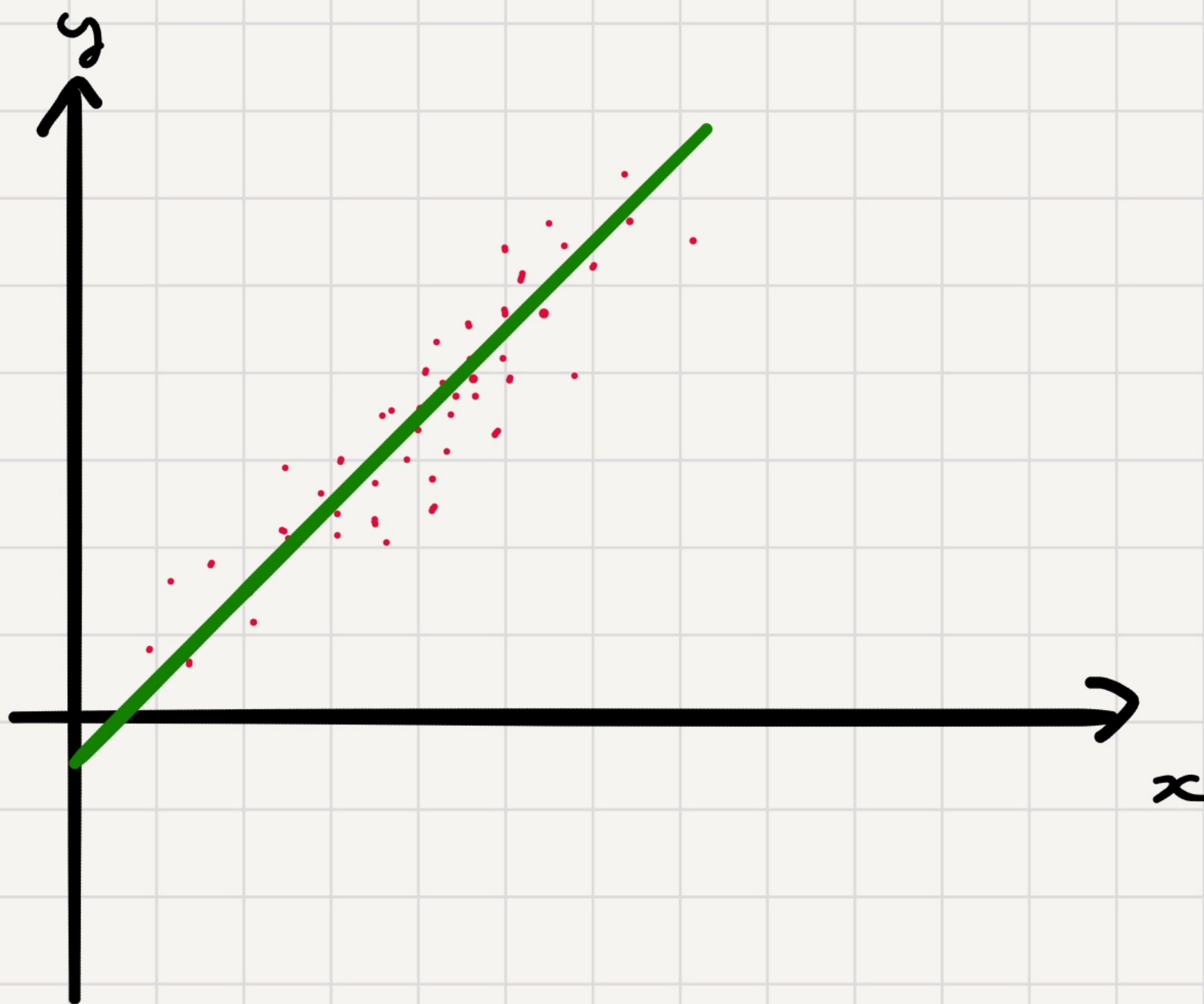
$$= E[(Y - EY)^2] - 2 \frac{\text{Cov}(X, Y)}{\text{Var}(X)} E[(Y - EY)(x - EX)]$$

$$+ \left(\frac{\text{Cov}(X, Y)}{\text{Var}(X)}\right)^2 E[(x - EX)^2]$$

$$= \text{Var}(Y) - \frac{\text{Cov}(X, Y)^2}{\text{Var}(X)}$$

$$= \text{Var}(Y) - \text{Var}(Y) \cdot \frac{\text{Cov}(X, Y)^2}{\text{Var}(X) \text{Var}(Y)} = \text{Var}(Y) \cdot (1 - \text{Corr}(X, Y)^2)$$

$$l(x) = \mathbb{E}Y + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} (x - \mathbb{E}X) \quad \left. \vphantom{\frac{\text{Cov}(X, Y)}{\text{Var}(X)}}} \right\}$$



The line goes through  $(\mathbb{E}X, \mathbb{E}Y)$   
The slope of the line is  $\frac{\text{Cov}(X, Y)}{\text{Var}(X)}$ . }

Theorem: The best linear predictor of  $Y$  as a function of  $X$  is  $EY + (x - EX) \cdot \frac{\text{cov}(X, Y)}{\text{var}(X)}$ .

Corollary: The minimum  $E[(Y - l(x))^2] = \text{Var}(Y) \cdot (1 - \text{Corr}(X, Y)^2)$ .

This is what we mean by

" $X$  explains 80% of the variance of  $Y$ "

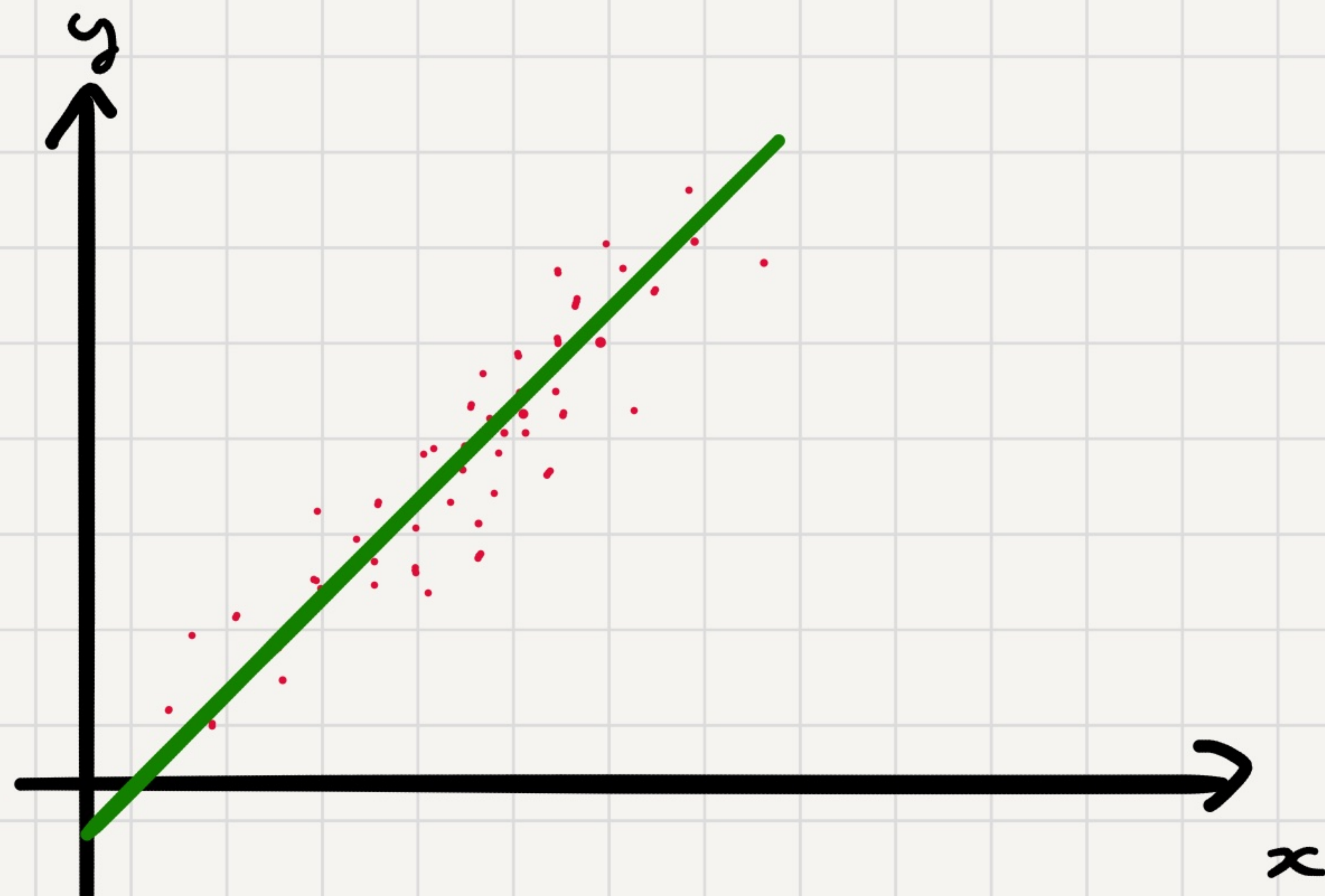


$$\text{corr}(X, Y)^2 = 0.8.$$

## Minimizing Given Data

Suppose you get samples  $(x_1, y_1), \dots, (x_n, y_n)$  from the joint distribution, and you want to minimize

$$\sum_{i=1}^n (y_i - \alpha x_i + \beta)^2$$



## Minimizing Given Data

Suppose you get samples  $(x_1, y_1), \dots, (x_n, y_n)$  from the joint distribution, and you want to minimize

$$\sum_{i=1}^n (y_i - \alpha x_i + \beta)^2$$

You'll get:

$$\bar{x} = \frac{1}{n} \cdot \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \cdot \sum_{i=1}^n y_i$$

$$\alpha = \frac{\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \cdot \bar{y}}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

(estimate to  $\frac{\text{Cov}(X, Y)}{\text{Var}(X)}$ )

$$\beta = \bar{y} - \alpha \cdot \bar{x}$$

(estimate to  $EY - \alpha EX$ )