



CS135

Introduction to Machine Learning

Lecture 4: Regression

Linear Regression

Two main objectives:

1. Establish relationship between 2 variables
 - Positive if ↑ in 1 causes ↑ in other
 - Negative if ↑ in 1 causes ↓ in other*e.g., Income & Spending*

2. Forecast new observations
Use relationship to predict unobserved events
e.g., Sales for next quarter

Linear Regression

Variable Roles:

- Dependent Variable
 - Value to explain or forecast
 - Depends on something else
 - Denoted as y
- Independent Variable
 - Variable to explain the other
 - Values independent
 - Denoted as x

Linear Regression

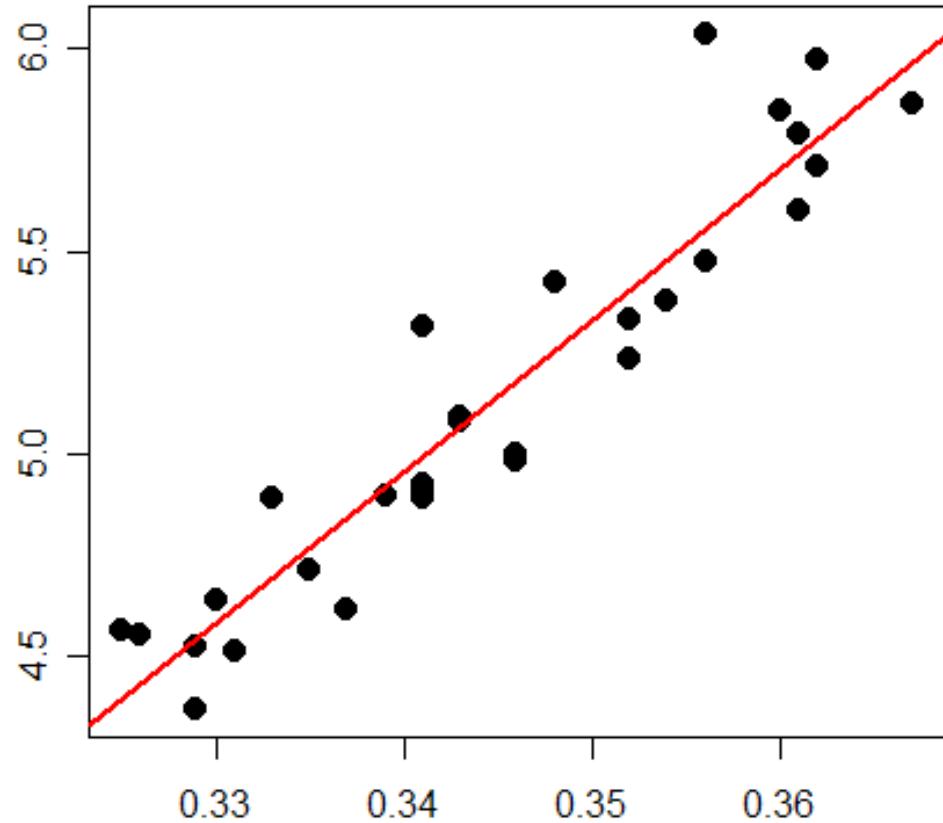
Linear Equation

$$y = a + bx$$

$$y = mx + b$$

$$y = \beta_0 + \beta_1 x$$

- β_0 is the intercept
- β_1 is the slope
- We will find these values via least squares



The Least Squares Criterion

The formulas for b_0 and b_1 that minimize the least squares criterion are:

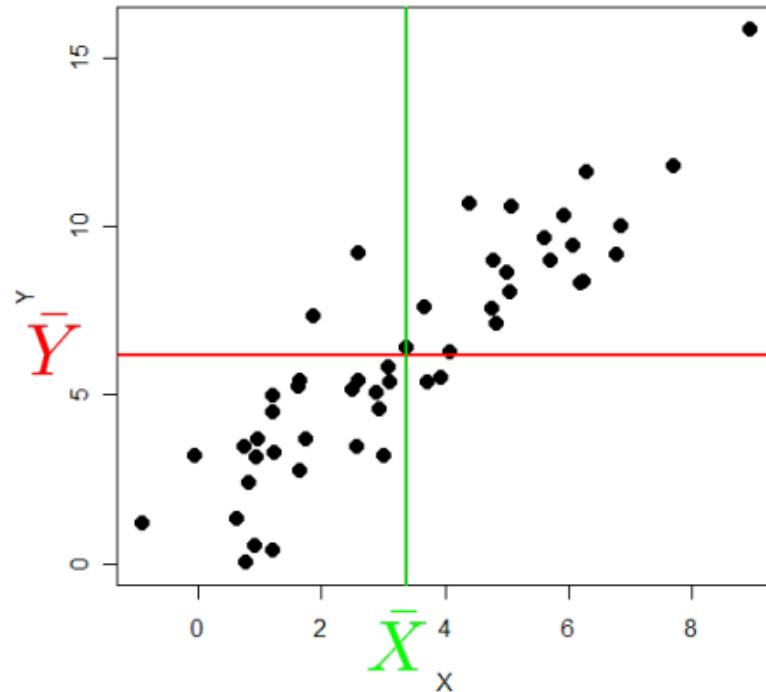
$$b_1 = \text{corr}(X, Y) \times \frac{s_Y}{s_X} \quad b_0 = \bar{Y} - b_1 \bar{X}$$

where,

$$s_Y = \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2} \quad \text{and} \quad s_X = \sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}$$

Correlation and Covariance

Measure the *direction* and *strength* of the linear relationship between variables Y and X



$$\text{Cov}(Y, X) = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{n - 1}$$

Correlation and Covariance

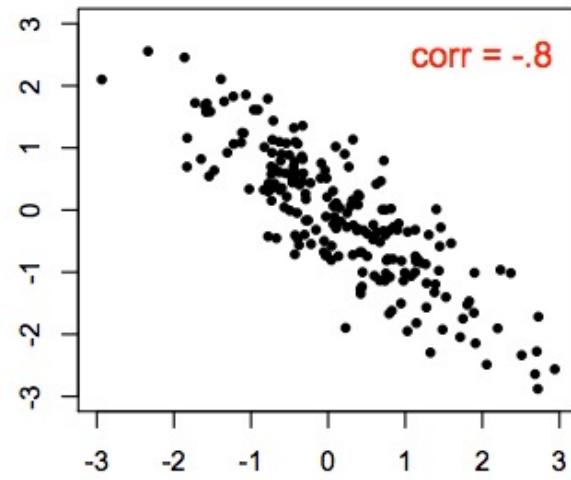
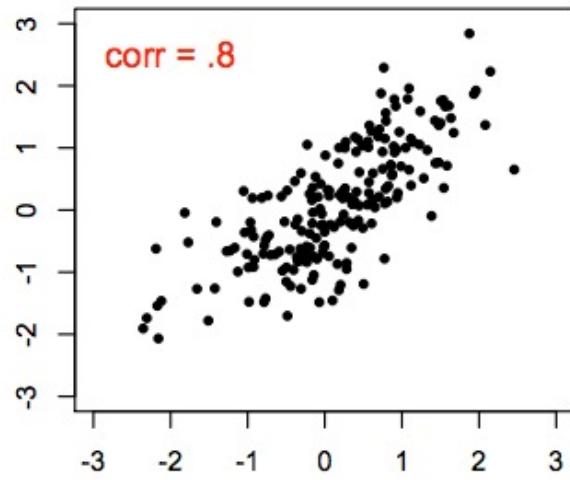
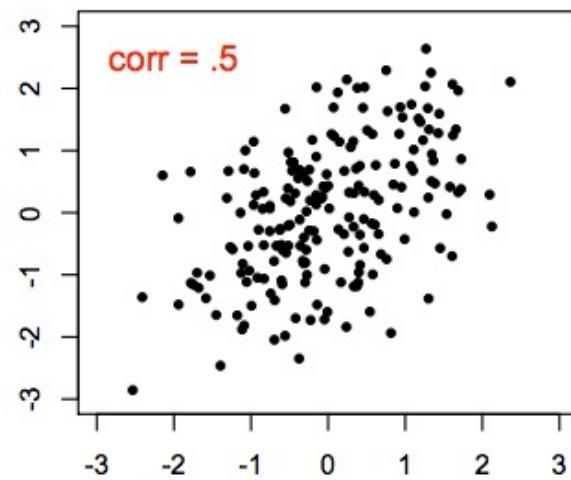
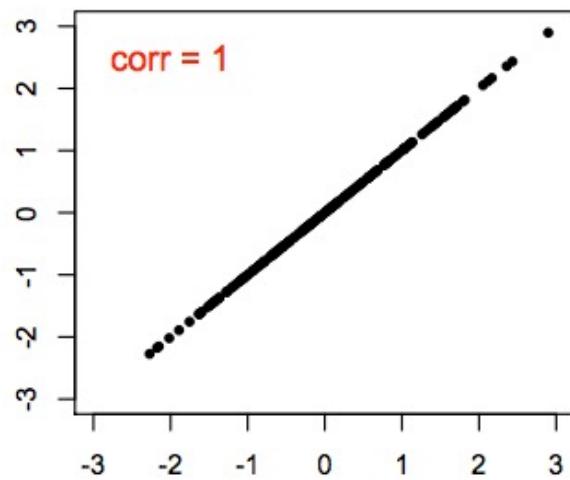
Correlation is the standardized covariance:

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}} = \frac{\text{cov}(X, Y)}{\text{sd}(X)\text{sd}(Y)}$$

The correlation is scale invariant and the units of measurement don't matter: It is always true that $-1 \leq \text{corr}(X, Y) \leq 1$.

This gives the direction (- or +) and strength ($0 \rightarrow 1$) of the linear relationship between X and Y .

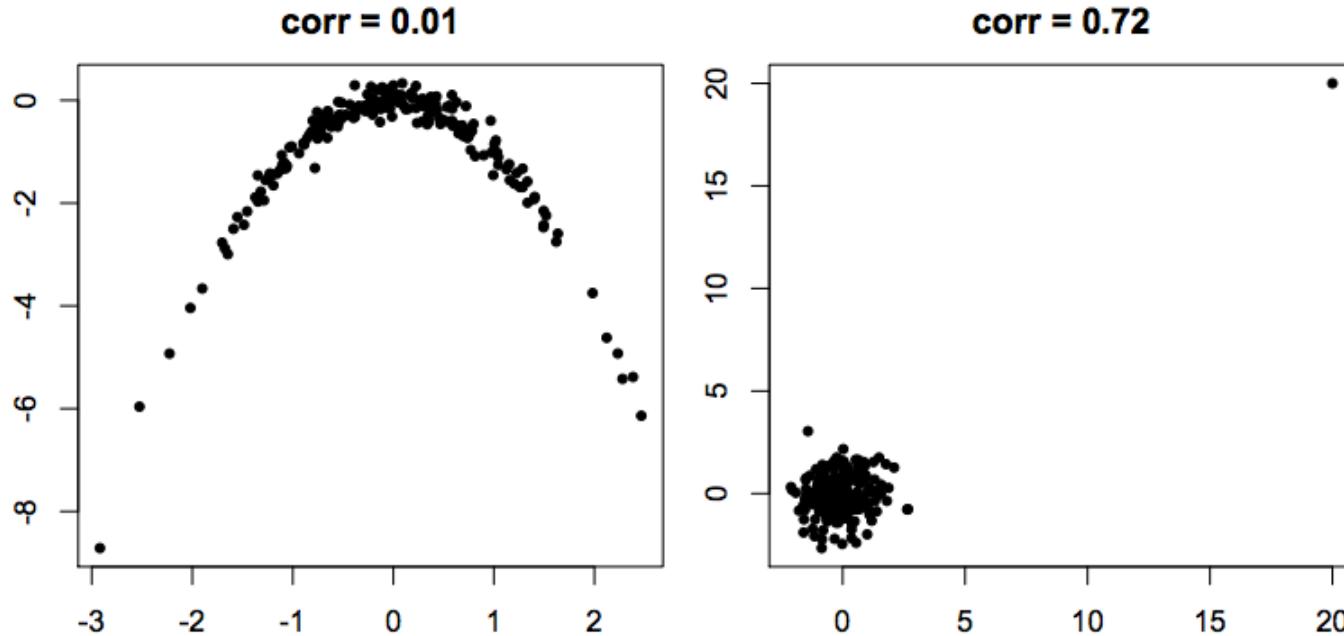
Correlation



Correlation

Only measures **linear** relationships:

$\text{corr}(X, Y) = 0$ does not mean the variables are not related!



Also be careful with influential observations.

Back To Least Squares

1. Intercept:

$$b_0 = \bar{Y} - b_1 \bar{X} \Rightarrow \bar{Y} = b_0 + b_1 \bar{X}$$

- ▶ The point (\bar{X}, \bar{Y}) is on the regression line!
- ▶ Least squares finds the point of means and rotate the line through that point until getting the “right” slope

2. Slope:

$$b_1 = \text{corr}(X, Y) \times \frac{s_Y}{s_X}$$

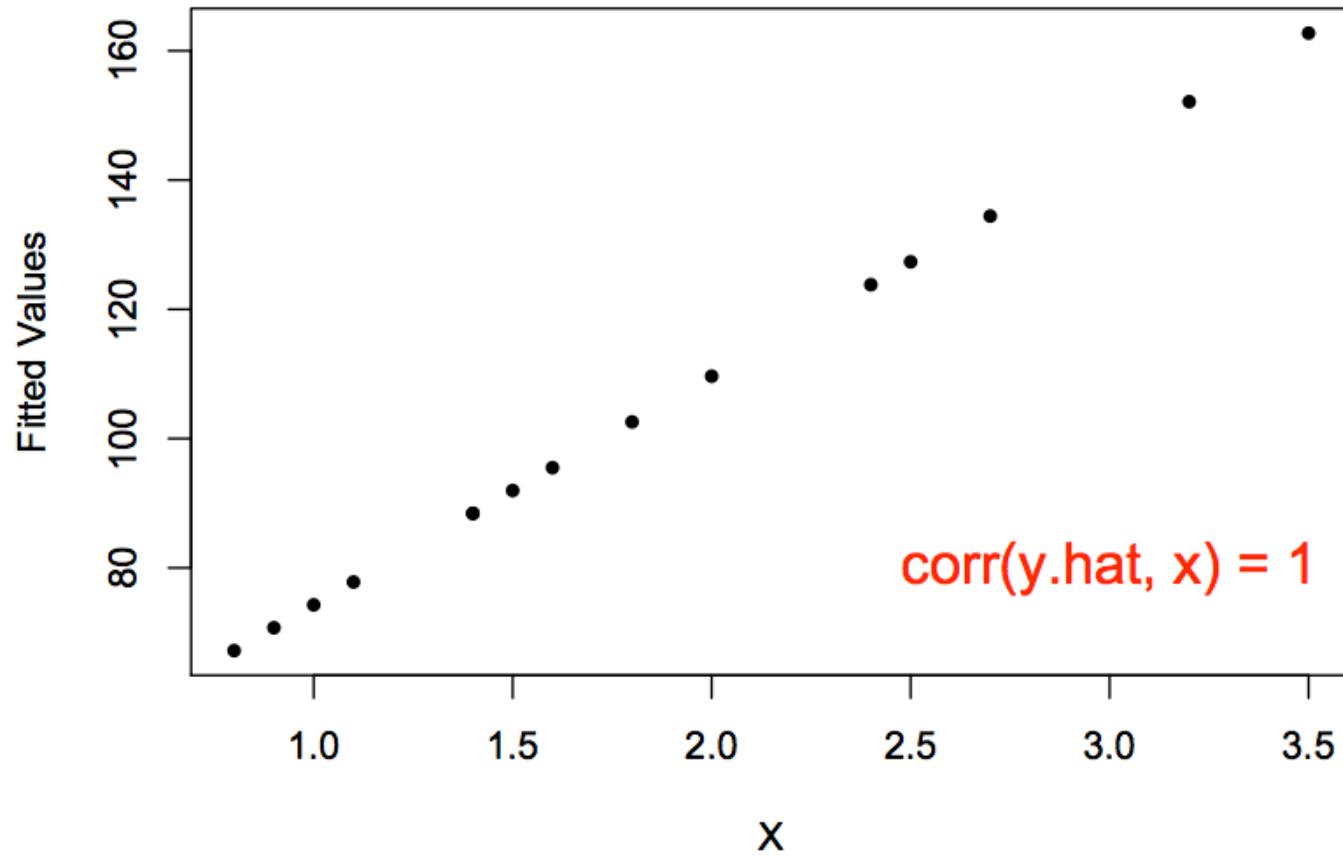
- ▶ So, the right slope is the *correlation coefficient* times a *scaling factor* that ensures the proper units for b_1

More on Least Squares

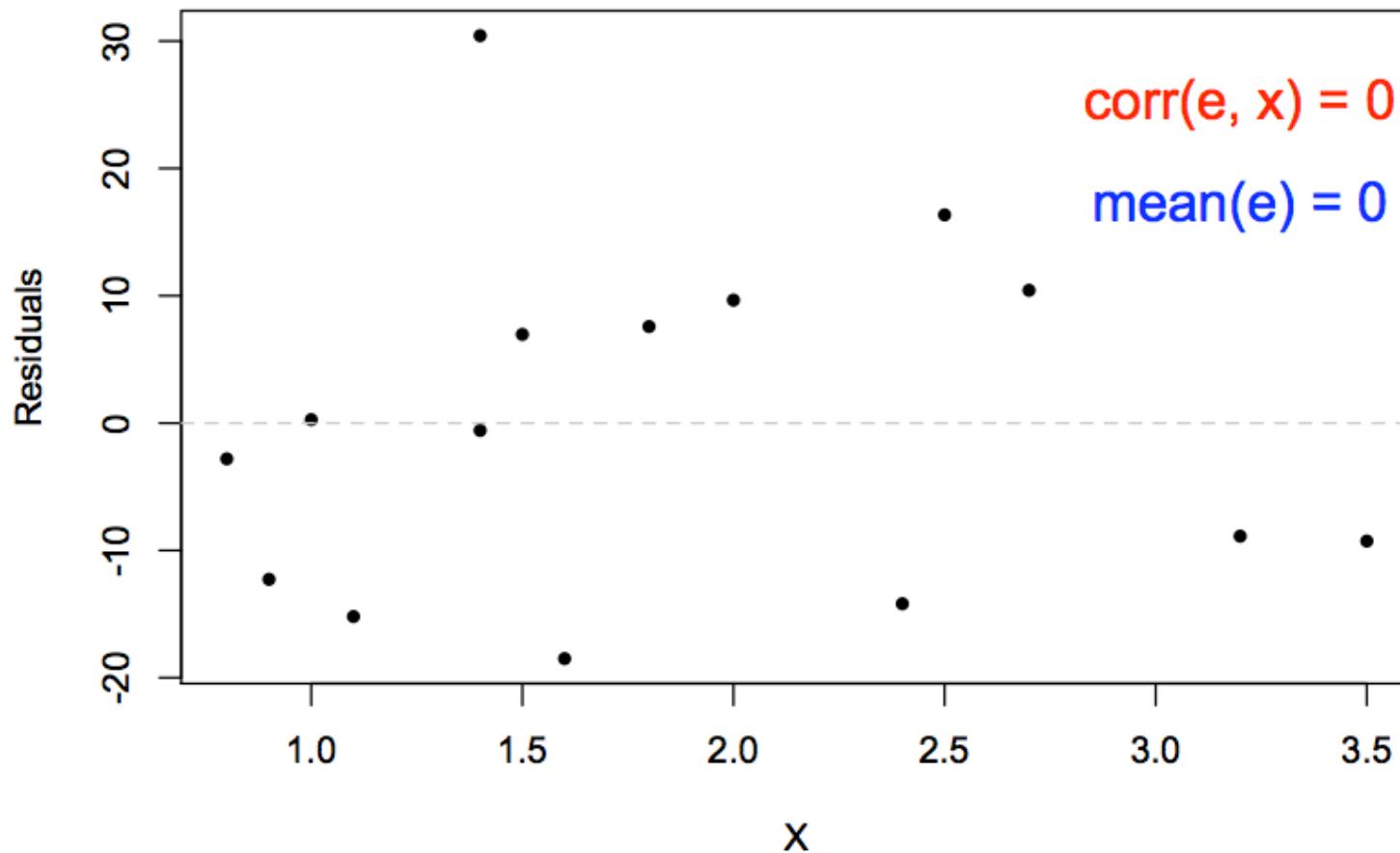
From now on, terms “fitted values” (\hat{Y}_i) and “residuals” (e_i) refer to those obtained from the least squares line.

The fitted values and residuals have some special properties. Lets look at the housing data analysis to figure out what these properties are...

The Fitted Values & X



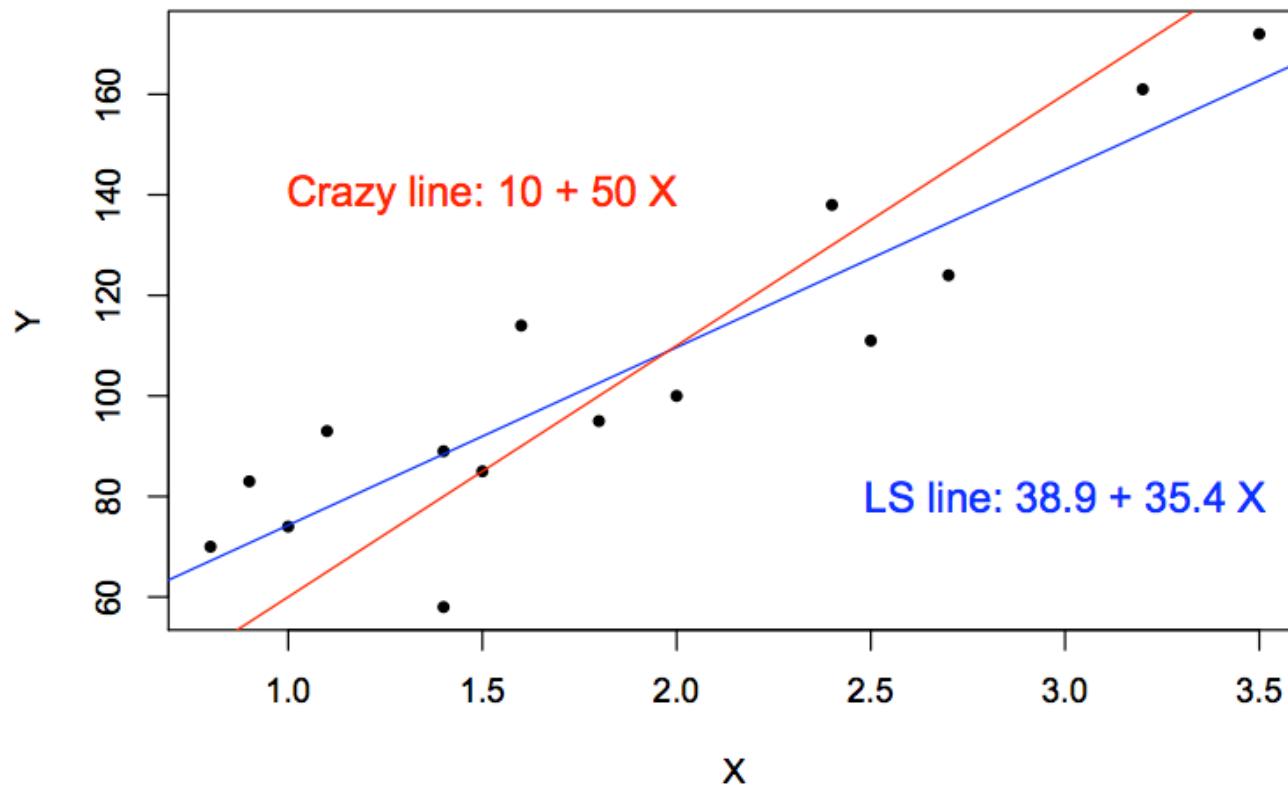
The Residuals of X



Why?

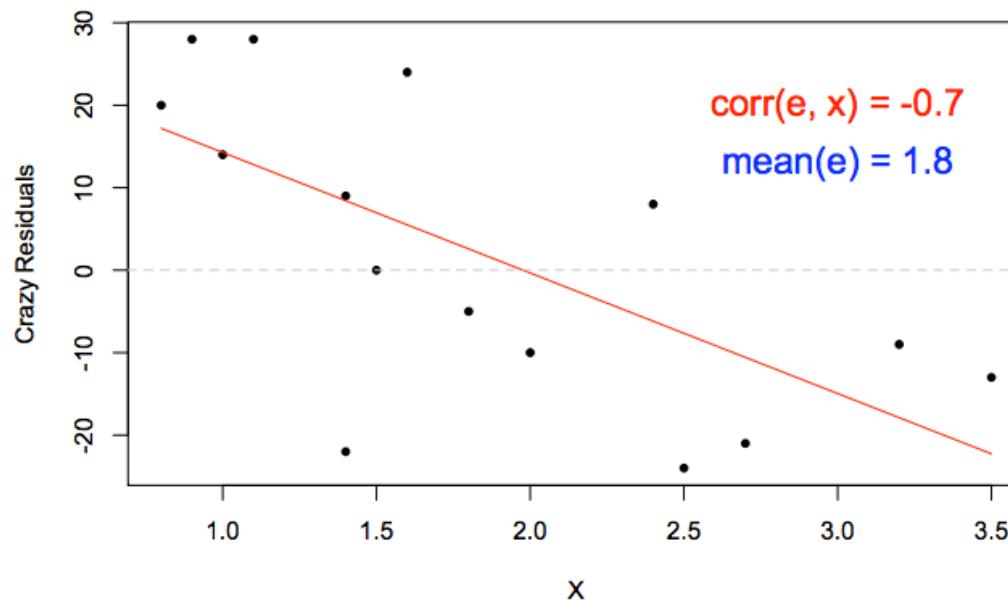
What is the intuition for the relationship between \hat{Y} and e and X ?

Lets consider some “crazy” alternative line:



Fitted Values and Residuals

This is a bad fit! We are underestimating the value of small houses and overestimating the value of big houses.



Clearly, we have left some predictive ability on the table!

Fitted Values and Residuals

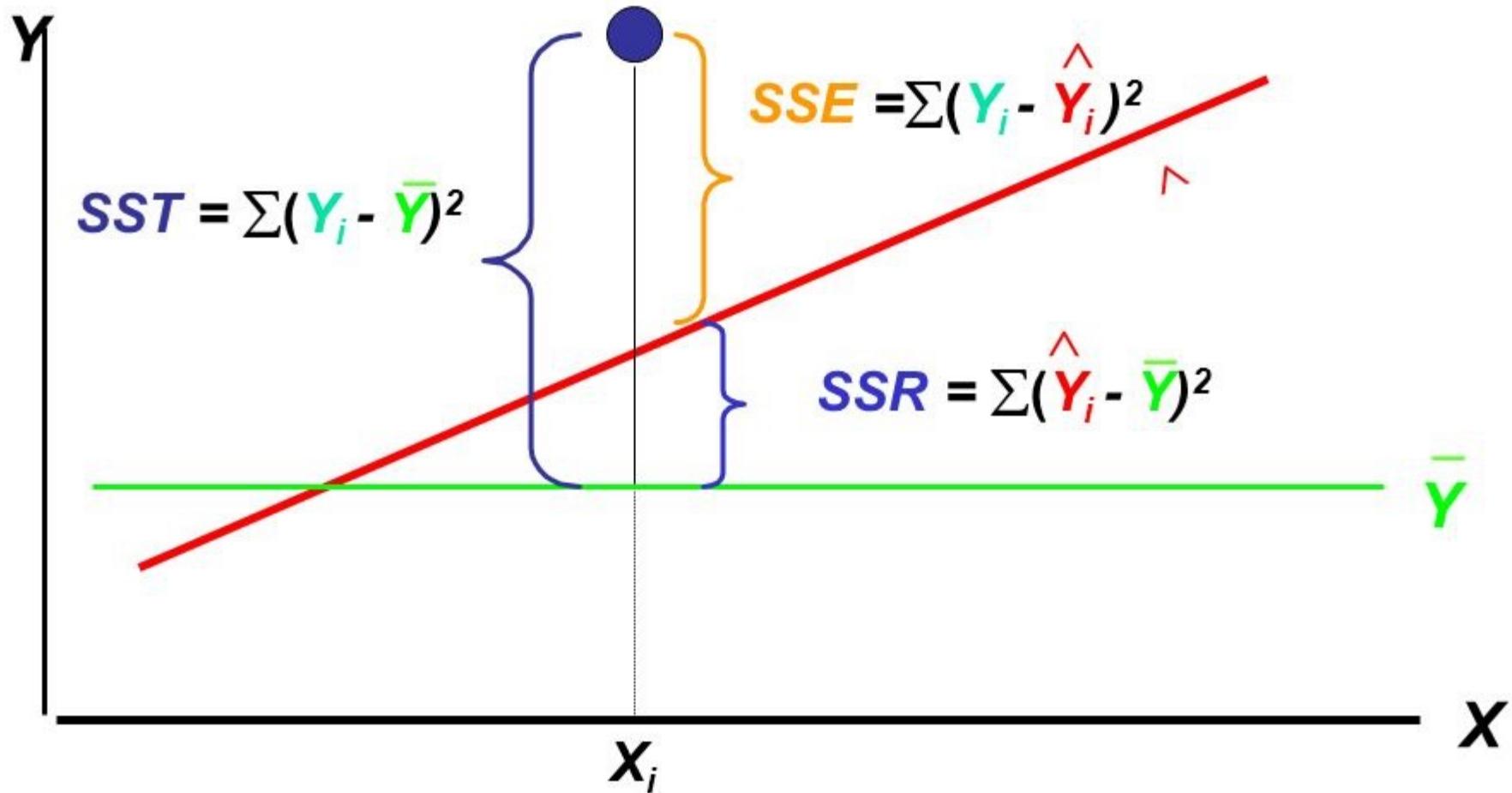
As long as the correlation between e and X is non-zero, we could always adjust our prediction rule to do better.

We need to exploit all of the predictive power in the X values and put this into \hat{Y} , leaving no “ X ness” in the residuals.

In Summary: $Y = \hat{Y} + e$ where:

- ▶ \hat{Y} is “made from X ”; $\text{corr}(X, \hat{Y}) = 1$.
- ▶ e is unrelated to X ; $\text{corr}(X, e) = 0$.

Measure Variation: The Sum of Squares



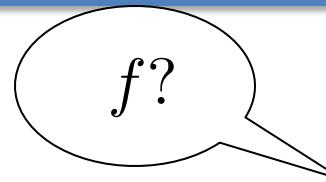
Regression

$x_i \in \mathbb{R}^d$ Gender Weight Age

Blood Pressure $y_i \in \mathbb{R}$

	Male	40K	19ys	3.1	
	Female	55K	34ys	1.2	
	⋮			⋮	
	Female	90K	24ys	2.5	

What for?



$$y_i \approx f(x_i), \quad i = 1, \dots, n$$

for some $f : \mathbb{R}^d \rightarrow \mathbb{R}$

Prediction: If $x = \boxed{\text{Female} \ 70K \ 20\text{ys}}$ then $y = ?$

Correlation: If Weight then $y \uparrow$

Linear Regression

$x_i \in \mathbb{R}^d$

	Gender	Weight	Age	Blood Pressure	$y_i \in \mathbb{R}$
n			40K	19ys	3.1
			55K	34ys	1.2
	⋮				
			90K	24ys	2.5

β ?



$$y_i \approx f(x_i) = \beta^\top x_i = \sum_{k=1}^d \beta_k x_{ik}$$

Assumption: There exists $\beta \in \mathbb{R}^d$ such that:

$$y_i = \langle \beta, x_i \rangle + \varepsilon_i, \quad i = 1, \dots, n$$

where ε_i are i.i.d., $\mathbb{E}[\varepsilon_i] = 0$, $\mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$

Linear Regression

$x_i \in \mathbb{R}^d$ Gender Weight Age

Blood Pressure $y_i \in \mathbb{R}$

	Gender	Weight	Age	Blood Pressure $y_i \in \mathbb{R}$
	♂	40K	19ys	
	♀	55K	34ys	
	⋮			
	♀	90K	24ys	3.1
				1.2
				2.5

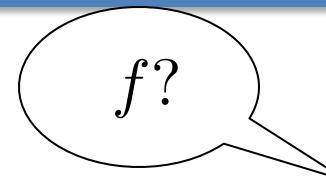
$$y_i \approx f(x_i), \quad i = 1, \dots, n$$

for some $f : \mathbb{R}^d \rightarrow \mathbb{R}$

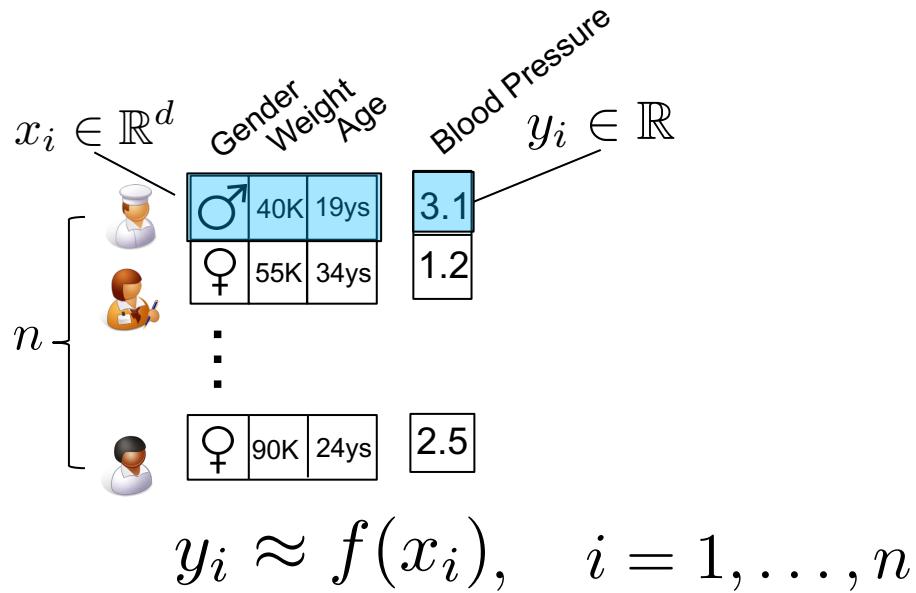
$$y = \underline{f(X)} + \varepsilon$$

Depends on statistical method

$$f(X) = \beta_0 + \beta_1 X$$



Regression: Terminology



- ❑ $x_i \in \mathbb{R}^d$: features, independent variables, covariates, inputs,...
- ❑ $y_i \in \mathbb{R}$: label, dependent variable, outcome, response, output,...

Regression

$x_i \in \mathbb{R}^d$

	Gender	Weight	Age
	Male	40K	19ys
	Female	55K	34ys
⋮			
	Female	90K	24ys

$y_i \in \mathbb{R}$

3.1	Blood Pressure
1.2	

$f?$



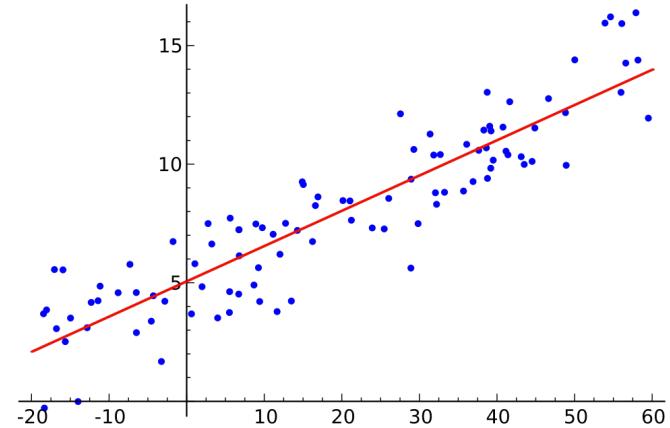
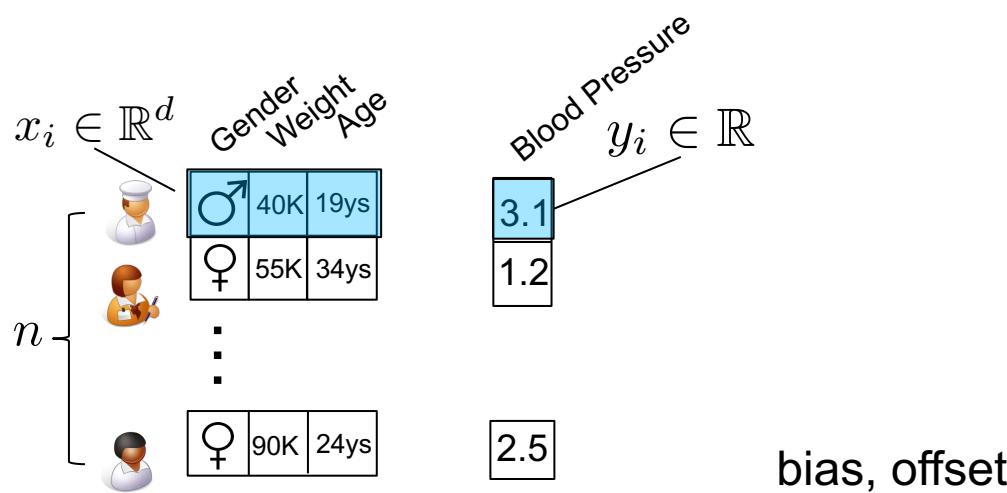
random "noise" variables

$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n.$$

where ε_i are **independent and identically distributed** (i.i.d), and

$$\mathbb{E}[\varepsilon_i] = 0 \quad \mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$$

Affine Can Be Written as Linear

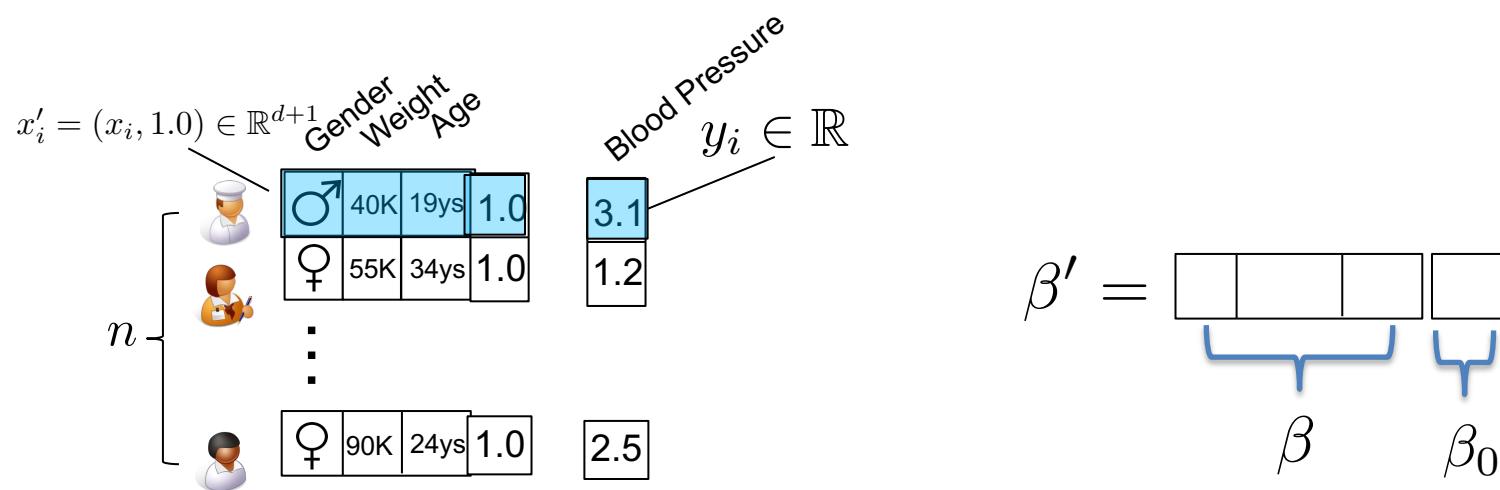


$$f(x_i) = \beta^\top x_i + \beta_0, \text{ where } \beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}$$

$$= \beta'^\top x'_i, \quad \text{where } \beta' = (\beta, \beta_0) \in \mathbb{R}^{d+1}$$

$$x'_i = (x_i, 1.0) \in \mathbb{R}^{d+1}$$

Affine Can Be Written as Linear

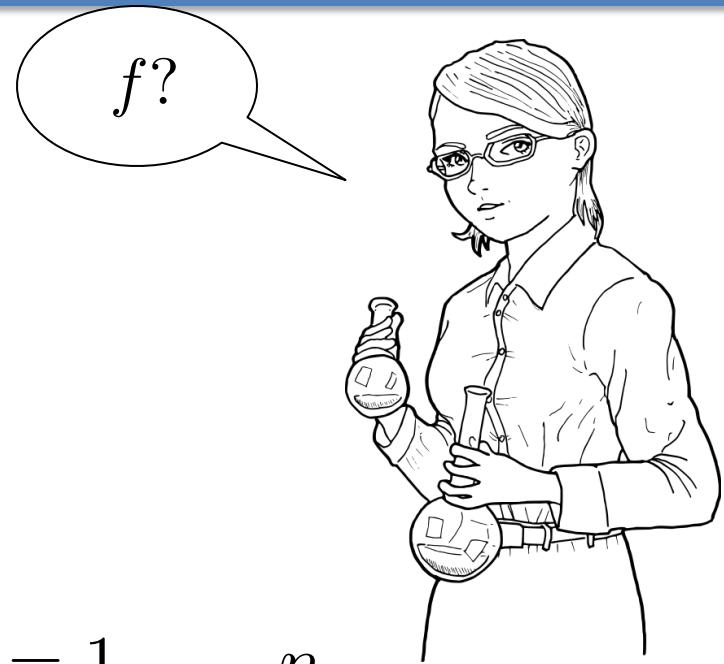
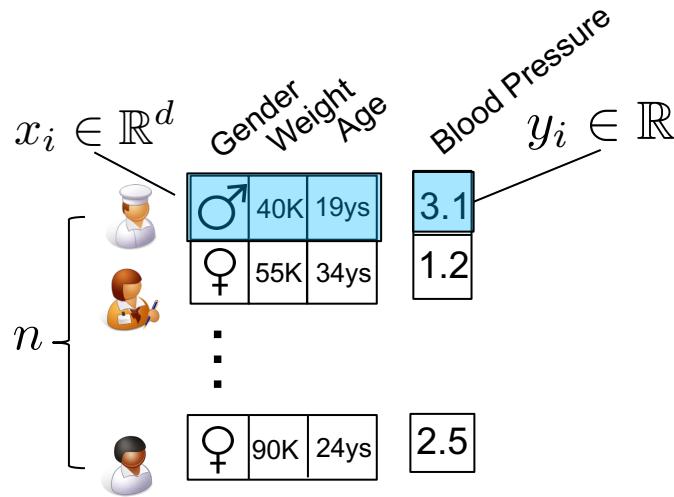


$$f(x_i) = \beta^\top x_i + \beta_0, \text{ where } \beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}$$

$$= \beta'^\top x'_i, \quad \text{where } \beta' = (\beta, \beta_0) \in \mathbb{R}^{d+1}$$

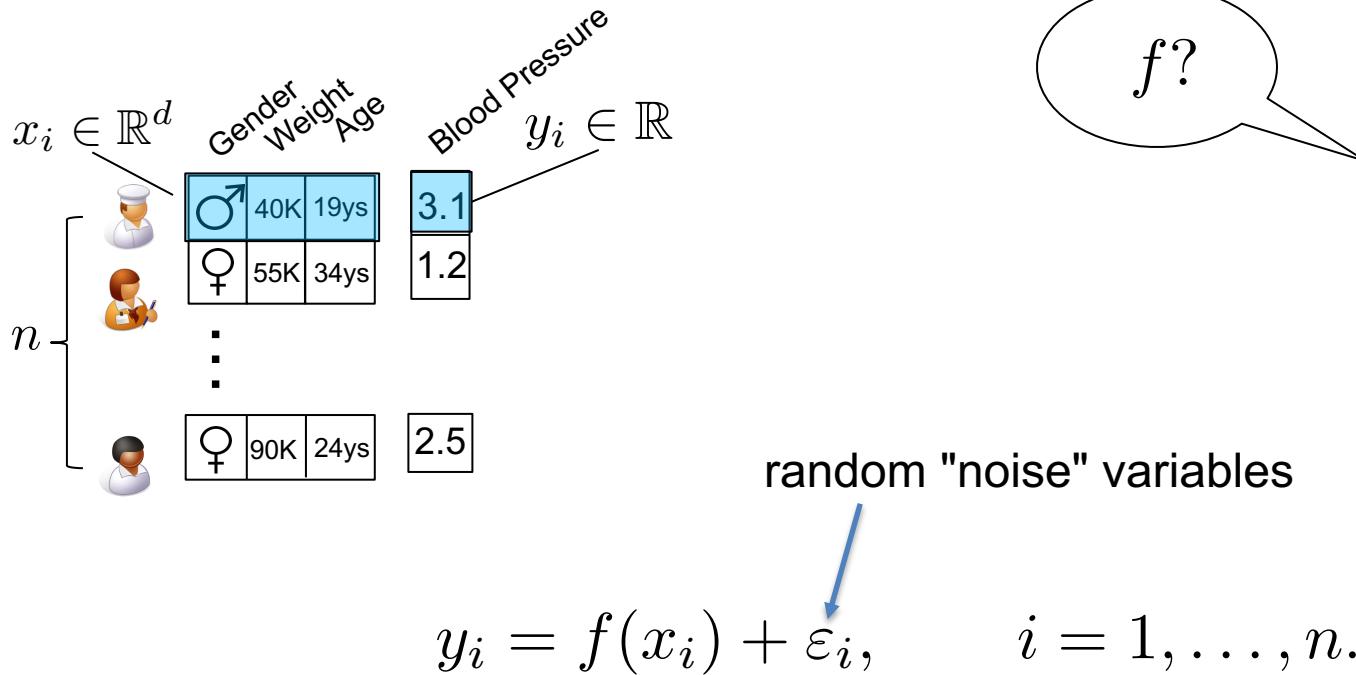
$$x'_i = (x_i, 1.0) \in \mathbb{R}^{d+1}$$

Regression: Noiseless Setting



$$y_i = f(x_i), \quad i = 1, \dots, n$$

Regression: Noisy setting



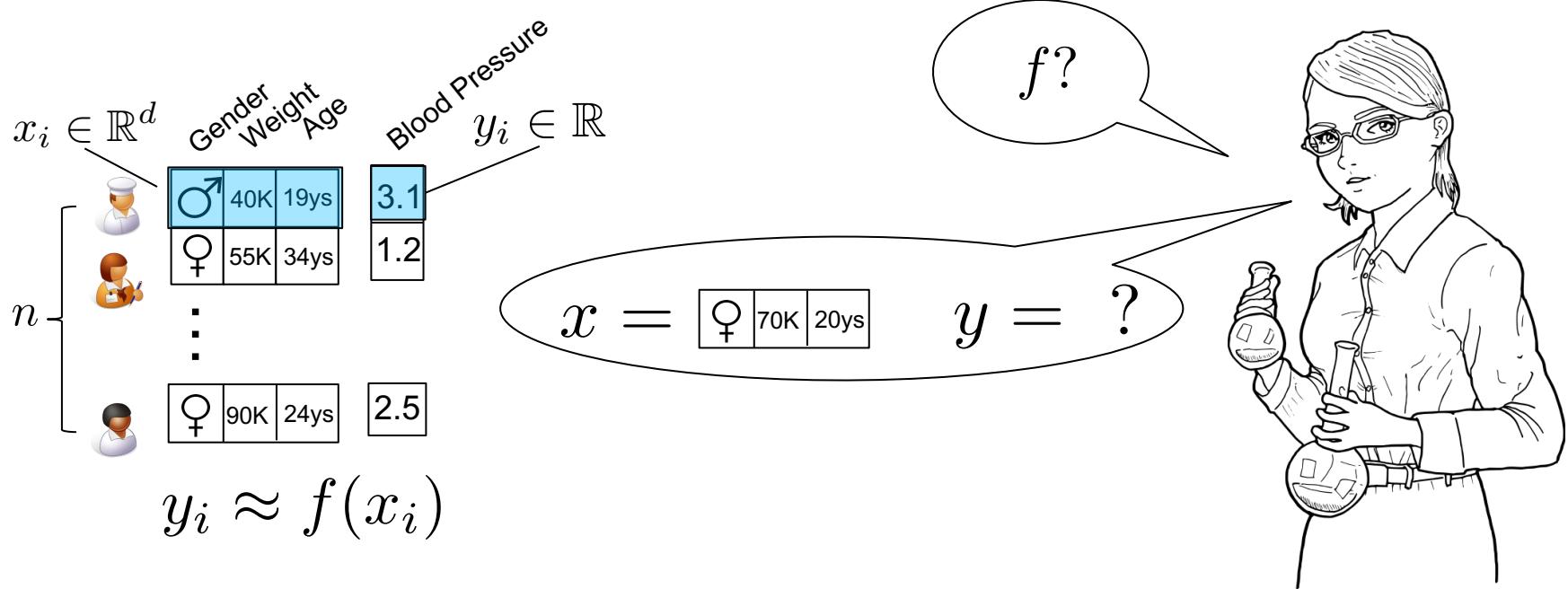
where ε_i are **independent and identically distributed** (i.i.d), and

$$\mathbb{E}[\varepsilon_i] = 0 \quad \mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$$

Note: This implies that $y_i, i = 1, \dots, n$, are **independent** random variables, where

$$\mathbb{E}[y_i] = f(x_i) \quad \text{Var}[y_i] = \mathbb{E} [(y_i - \mathbb{E}[y_i])^2] = \sigma^2$$

How would you solve this problem?



...you need to start making some assumptions on f !

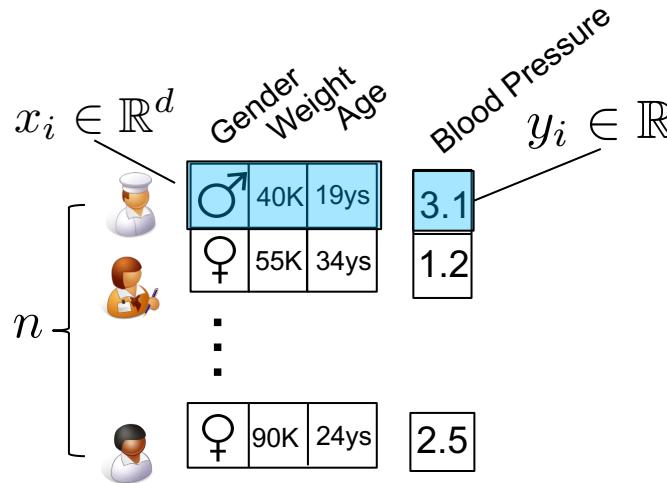
Assumption: Continuity!

□ **Assumption:** Function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is continuous

If $\lim_{k \rightarrow \infty} x_k = x$ then $\lim_{k \rightarrow \infty} f(x_k) = f(x)$

□ Values of f at points near x tell you something about $f(x)$!

Operations that Preserve Linearity



- Affine in \mathbb{R}^d is linear in \mathbb{R}^{d+1}
- Linear transforms on features (i.e., rescaling):
 - E.g., from kilograms to pounds
- Affine transforms in features (i.e., rescaling and shifting):
 - E.g., from F° to C° .

Linear Regression (return)

- Linear Regression Estimate

$$\hat{y} = \widehat{f(X)}$$

$$\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 X$$

- Linear Regression Actual Response

$$y = f(X) + \epsilon$$

Linear Regression

- Residual Error

$$e = y - \hat{y}$$

$$\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 X$$

- Do not confuse residual error with irreproducible error

$$e = \beta_0 + \beta_1 X + \epsilon - (\widehat{\beta}_0 + \widehat{\beta}_1 X)$$

$$e = (\beta_0 - \widehat{\beta}_0) + (\beta_1 - \widehat{\beta}_1)X + \epsilon$$

Linear Regression

$$\{(X_1, y_1), (X_2, y_2), \dots, (X_n, y_n)\}$$

- Residual Error of i^{th} sample

$$e_i = y_i - \hat{y}_i$$

- Sum of Squares Residual (RSS)

$$RSS = \sum_{i=1}^n e_i^2$$

Linear Regression

$$\operatorname{argmin}_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \hat{y})^2$$

$$\operatorname{argmin}_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - (\widehat{\beta}_0 + \widehat{\beta}_1 X_i))^2$$

$$\operatorname{argmin}_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \widehat{\beta}_0 - \widehat{\beta}_1 X_i)^2$$

Linear Regression

$$\operatorname{argmin}_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \widehat{\beta}_0 - \widehat{\beta}_1 X_i)^2$$

□ Differentiating wrt β_0

$$2 \sum_{i=1}^n (y_i - \widehat{\beta}_0 - \widehat{\beta}_1 X_i) (-1) = 0$$

$$\sum_{i=1}^n (y_i - \widehat{\beta}_0 - \widehat{\beta}_1 X_i) = 0$$

$$\widehat{\beta}_0 = \frac{\sum_{i=1}^n (y_i - \widehat{\beta}_1 X_i)}{n}$$

$$\widehat{\beta}_0 = \sum_{i=1}^n \left(\frac{y_i}{n} - \widehat{\beta}_1 \frac{X_i}{n} \right)$$

$$\widehat{\beta}_0 = \sum_{i=1}^n \frac{y_i}{n} - \widehat{\beta}_1 \sum_{i=1}^n \frac{X_i}{n}$$

$$\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \bar{x}$$

Linear Regression

$$\operatorname{argmin}_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \widehat{\beta}_0 - \widehat{\beta}_1 X_i)^2$$

□ Differentiating wrt β_1

$$2 \sum_{i=1}^n (y_i - \widehat{\beta}_0 - \widehat{\beta}_1 X_i) (-X_i) = 0$$

$$\sum_{i=1}^n (y_i - \widehat{\beta}_0 - \widehat{\beta}_1 X_i) (X_i) = 0$$

$$\sum_{i=1}^n X_i y_i - \widehat{\beta}_0 \sum_{i=1}^n X_i - \widehat{\beta}_1 \sum_{i=1}^n X_i^2 = 0$$

$$\sum_{i=1}^n X_i y_i - \left(\sum_{i=1}^n \frac{y_i}{n} - \widehat{\beta}_1 \sum_{i=1}^n \frac{X_i}{n} \right) \sum_{i=1}^n X_i - \widehat{\beta}_1 \sum_{i=1}^n X_i^2 = 0$$

$$\sum_{i=1}^n X_i y_i - \sum_{i=1}^n \frac{y_i}{n} \sum_{i=1}^n X_i + \widehat{\beta}_1 \sum_{i=1}^n \frac{X_i}{n} \sum_{i=1}^n X_i - \widehat{\beta}_1 \sum_{i=1}^n X_i^2 = 0$$

$$\sum_{i=1}^n X_i y_i - \frac{1}{n} \sum_{i=1}^n y_i \sum_{i=1}^n X_i - \widehat{\beta}_1 \times \frac{1}{n} \sum_{i=1}^n X_i \sum_{i=1}^n X_i - \widehat{\beta}_1 \sum_{i=1}^n X_i^2 = 0$$

$$\widehat{\beta}_1 = \frac{\sum_{i=1}^n X_i y_i - \frac{1}{n} \sum_{i=1}^n y_i \sum_{i=1}^n X_i}{\sum_{i=1}^n X_i^2 - \frac{1}{n} (\sum_{i=1}^n X_i)^2}$$

$$\boxed{\widehat{\beta}_1 = \frac{\sum_{i=1}^n X_i y_i - n \bar{X} \bar{y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2}}$$

Linear Regression

$$\{(X_1, y_1), (X_2, y_2), \dots, (X_n, y_n)\}$$

Least Squares Criteria

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n X_i y_i - n \bar{X} \bar{y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2}$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

Car Sales

Number of Ads Run

Multiple Regression

- General Parametric Equation

$$\hat{y} = \underline{f(\mathbf{X}) + \epsilon}$$

Depends on statistical method

$$f(\mathbf{X}) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p$$

$$\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1 X_1 + \cdots + \widehat{\beta}_p X_p$$

- Number of operations = $n \times (p - 1)^2$, where n is the number of samples

Multiple Regression

$$\begin{array}{c} n \times 1 \\ n \times (p+1) \\ (p+1) \times 1 \\ n \times 1 \end{array} \quad \begin{array}{l} \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & X_{1,1} & X_{1,2} & \cdots & \cdots & X_{1,p} \\ 1 & X_{2,1} & X_{2,2} & \cdots & \cdots & \cdot \\ 1 & X_{3,1} & X_{3,2} & \cdots & \cdots & \cdot \\ 1 & \cdot & \cdot & \ddots & \ddots & \cdot \\ 1 & \cdot & \cdot & \ddots & \ddots & \cdot \\ 1 & X_{n,1} & X_{n,2} & \cdots & \cdots & X_{n,p} \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{bmatrix} \end{array}$$

$$y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

$$\hat{y} = \mathbf{X}\hat{\boldsymbol{\beta}}$$

Multiple Regression

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

$$\mathbf{e} = \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \vdots \\ \vdots \\ e_n \end{bmatrix} = \begin{bmatrix} y_1 - \widehat{y}_1 \\ y_2 - \widehat{y}_2 \\ y_3 - \widehat{y}_3 \\ \vdots \\ \vdots \\ y_n - \widehat{y}_n \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \widehat{y}_1 \\ \widehat{y}_2 \\ \widehat{y}_3 \\ \vdots \\ \vdots \\ \widehat{y}_n \end{bmatrix} = \mathbf{y} - \widehat{\mathbf{y}}$$

$$RSS = \sum_{i=1}^n e_i^2 \quad \longrightarrow \quad RSS = \mathbf{e}^T \mathbf{e}$$

Multiple Regression

$$RSS = e^T e$$

$$RSS = (y - \hat{y})^T (y - \hat{y})$$

$$RSS = (y - X\hat{\beta})^T (y - \widehat{X}\hat{\beta})$$

$$RSS = (y - \hat{\beta}^T X^T) (y - \widehat{X}\hat{\beta})$$

$$RSS = y^T y^T - y^T X \hat{\beta} - \hat{\beta} X^T y + \hat{\beta}^T X^T X \hat{\beta}$$

Multiple Regression

$$RSS = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}} \mathbf{X}^T \mathbf{y} + \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}}$$

Matrix Differentiation

$x = m \times 1$ matrix

$A = n \times m$ matrix; $A \perp x$

$$\mathbf{y} = A \rightarrow \frac{\delta \mathbf{y}}{\delta x} = \mathbf{0}$$

$$\mathbf{y} = Ax \rightarrow \frac{\delta \mathbf{y}}{\delta x} = A$$

$$\mathbf{y} = xA \rightarrow \frac{\delta \mathbf{y}}{\delta x} = A^T$$

$$\mathbf{y} = x^T Ax \rightarrow \frac{\delta \mathbf{y}}{\delta x} = 2x^T A$$

$$\frac{\delta(RSS)}{\delta \hat{\boldsymbol{\beta}}} = \frac{\delta(\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}} \mathbf{X}^T \mathbf{y} + \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}})}{\delta \hat{\boldsymbol{\beta}}} = 0$$

$$\frac{\delta(\mathbf{y}^T \mathbf{y})}{\delta \hat{\boldsymbol{\beta}}} - \frac{\delta(\mathbf{y}^T \mathbf{X} \hat{\boldsymbol{\beta}})}{\delta \hat{\boldsymbol{\beta}}} - \frac{\delta(\hat{\boldsymbol{\beta}} \mathbf{X}^T \mathbf{y})}{\delta \hat{\boldsymbol{\beta}}} + \frac{\delta(\hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}})}{\delta \hat{\boldsymbol{\beta}}} = 0$$

$$0 - \mathbf{y}^T \mathbf{X} - (\mathbf{X}^T \mathbf{y})^T + 2\hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} = 0$$

$$2\hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} = 2\mathbf{y}^T \mathbf{X} \longrightarrow \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} = \mathbf{y}^T \mathbf{X}$$

$$\hat{\boldsymbol{\beta}}^T = \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \longrightarrow \boxed{\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}$$

Multiple Regression

High Bias Problem

$$\{(X_1, y_1), (X_2, y_2), \dots, (X_n, y_n)\}$$

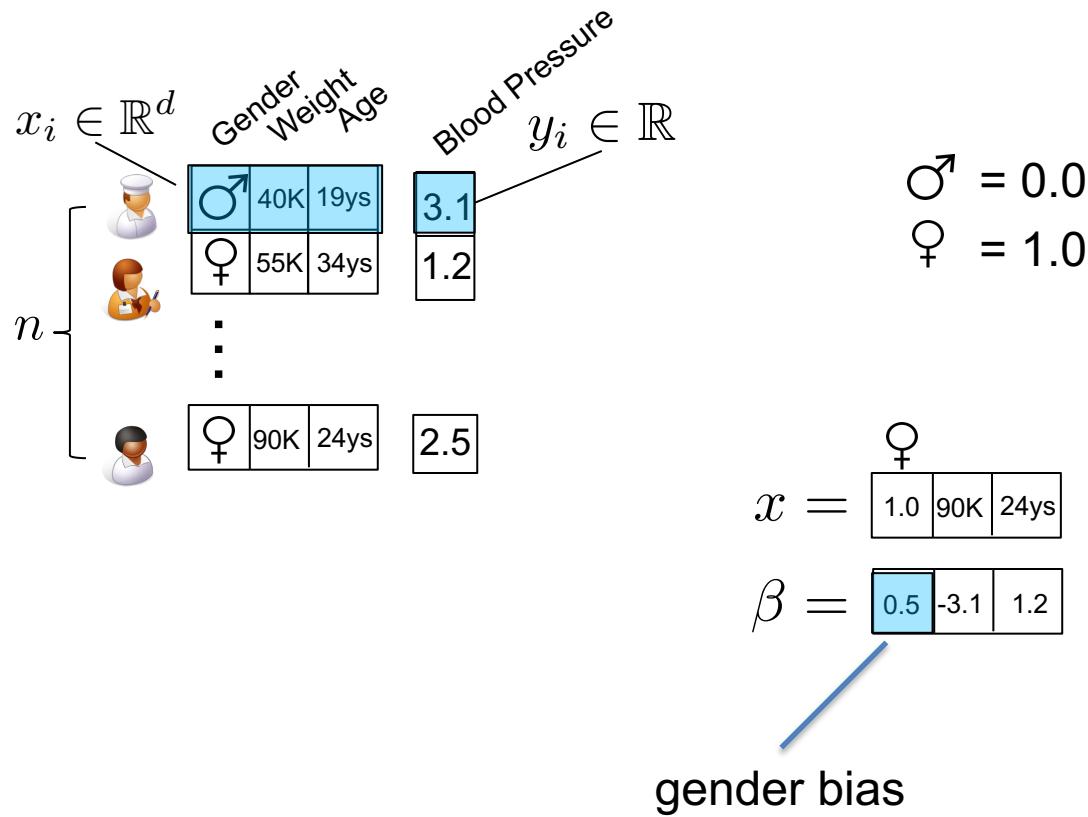
Least Squares Criteria

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

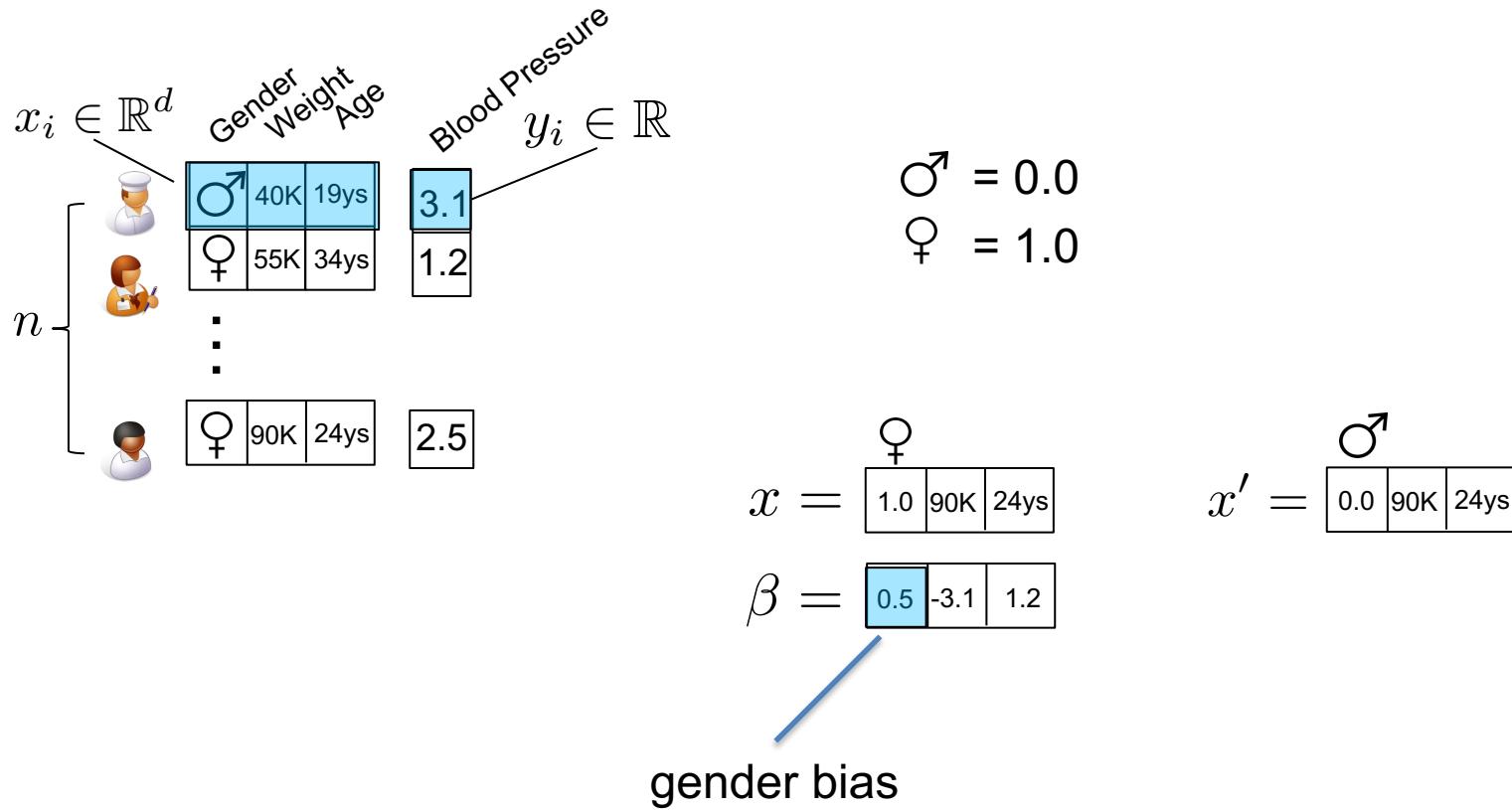
$$\hat{y} = \widehat{f(X)} = X \hat{\beta}$$

Car Sales Number of Ads Run
Number of Salesman
⋮
feature_p

Binary Features

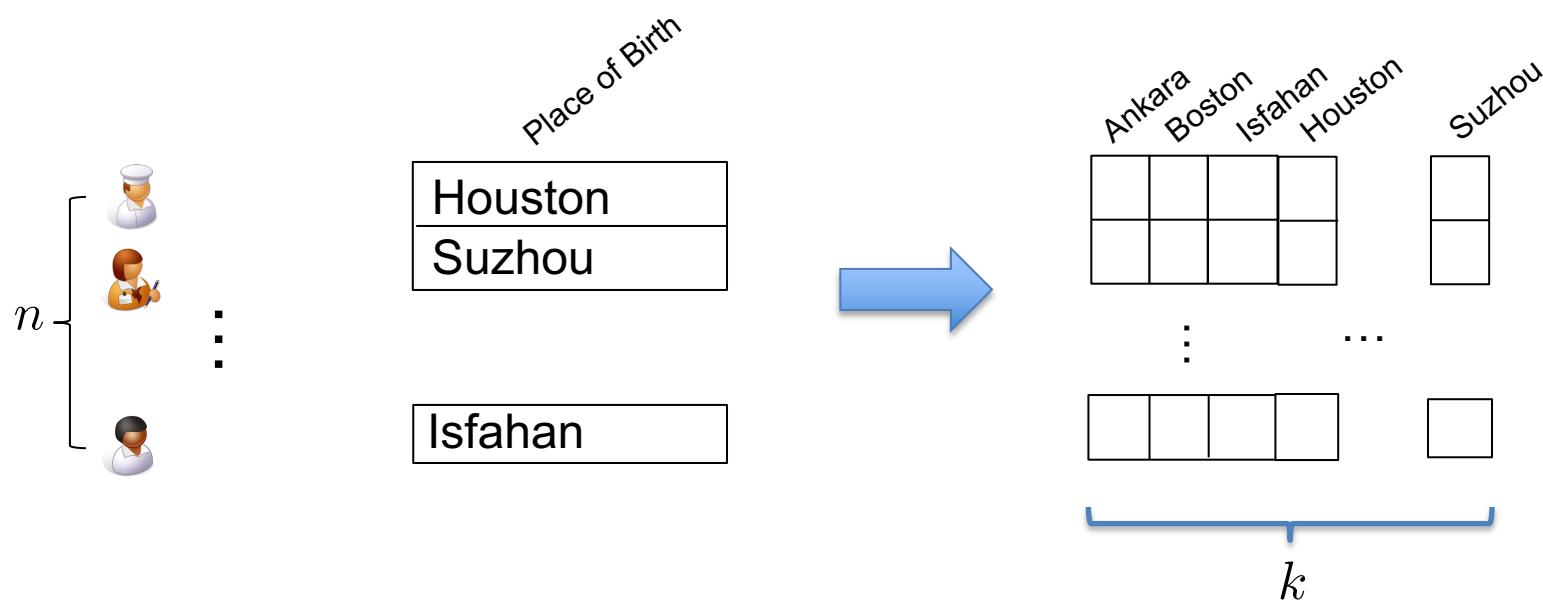


Binary Features



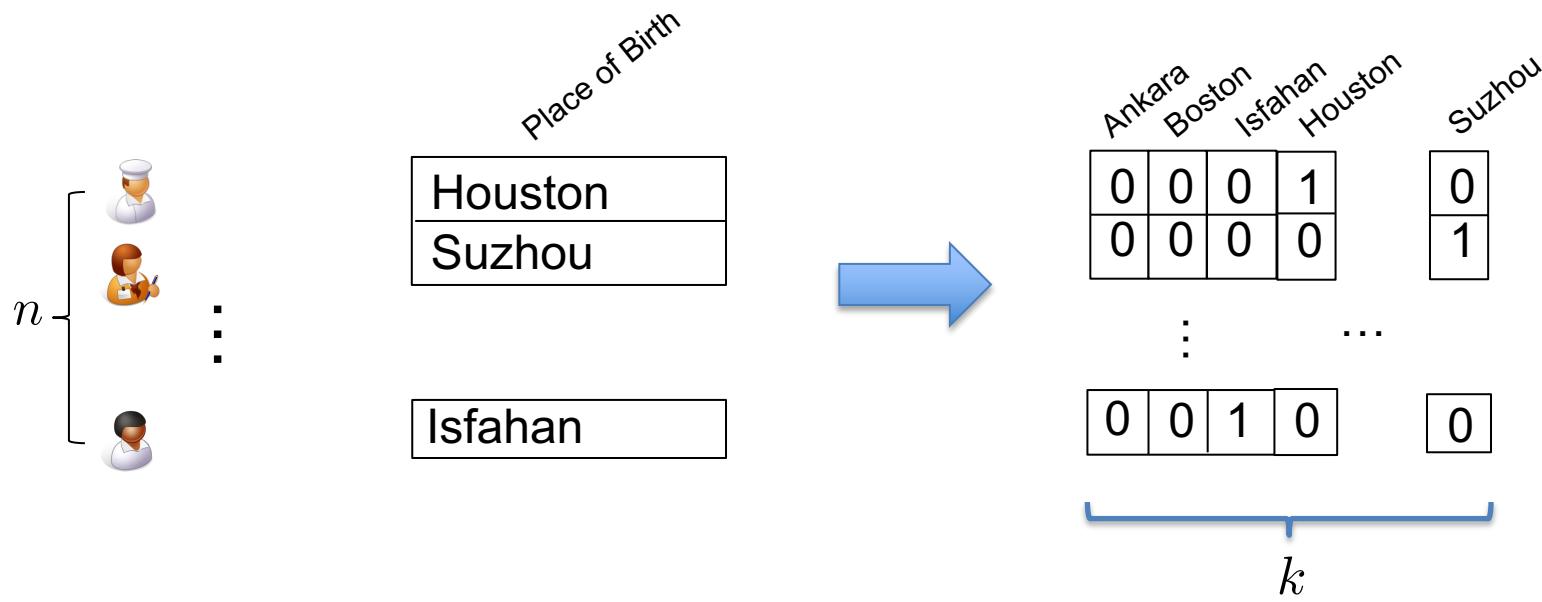
$$\mathbb{E}[y - y'] = \beta^\top x - \beta^\top x' = 0.5$$

Categorial Features: Binarization



Categories (cities) = k

Categorial Features: Binarization



$$\# \text{ Categories (cities)} = k$$

- ❑ Categorical features are very common: locations, genes, words in document
- ❑ Binarization leads to feature vectors that are **sparse**: most elements are 0!

Numeric Features May Be Categorical!!!!

n	ZIPCODE	Day of Week Examined	Blood Pressure
	02115	1	3.1
	02130	6	1.2
	02122	5	2.5

Mon: 1

Tue: 2

...

Sun: 7

Rule of thumb: if 2 does not mean "2 times" 1, treat it as categorical

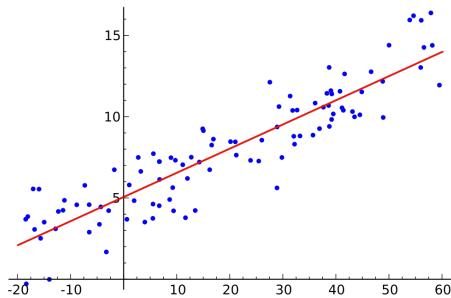
Least Squares Estimator (LSE)

$x_i \in \mathbb{R}^d$ Gender Weight Age

Blood Pressure $y_i \in \mathbb{R}$

	Male	40K	19ys	3.1
	Female	55K	34ys	1.2
	⋮			
	Female	90K	24ys	2.5

$$y_i \approx f(x_i) = \beta^\top x_i = \sum_{k=1}^d \beta_k x_{ik}$$



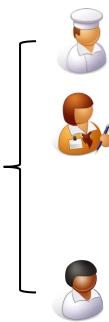
Estimate of β

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \langle \beta, x_i \rangle)^2$$

$\beta ?$



Least Squares Estimator (LSE)

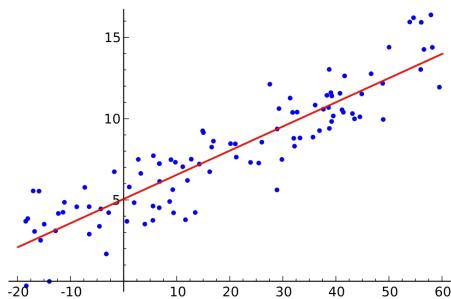


Gender
Weight
Age

Blood Pressure

$$X \in \mathbb{R}^{n \times d}$$
$$y \in \mathbb{R}^n$$

$$y_i \approx f(x_i) = \beta^\top x_i = \sum_{k=1}^d \beta_k x_{ik}$$



Why LSE?

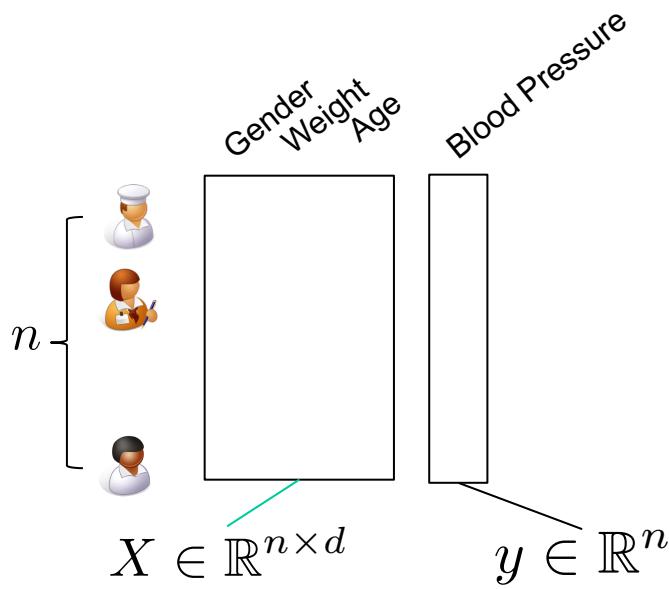
Estimate of β

$$\begin{aligned}\hat{\beta} &= \arg \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \langle \beta, x_i \rangle)^2 \\ &= \arg \min_{\beta \in \mathbb{R}^d} \|X\beta - y\|_2^2 \\ &= (X^T X)^{-1} X^T y\end{aligned}$$

$\beta ?$



Reason: If Noise is Gaussian, LSE is an MLE!



$$y_i = \beta^\top x_i + \varepsilon_i, \quad i = 1, \dots, n$$

$$\varepsilon_i \text{ i.i.d., } \mathbb{E}[\varepsilon_i] = 0, \mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$$

□ Suppose, in addition, that

$$\varepsilon_i \sim N(0, \sigma^2)$$

Then, the negative log-likelihood of the labels is:

$$\begin{aligned} -\log(P(y|\beta, X)) &= -\log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-(y_i - \beta^\top x_i)^2 / 2\sigma^2} \right) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^\top x_i)^2 + C \end{aligned}$$

Bias vs. Variance Trade-off

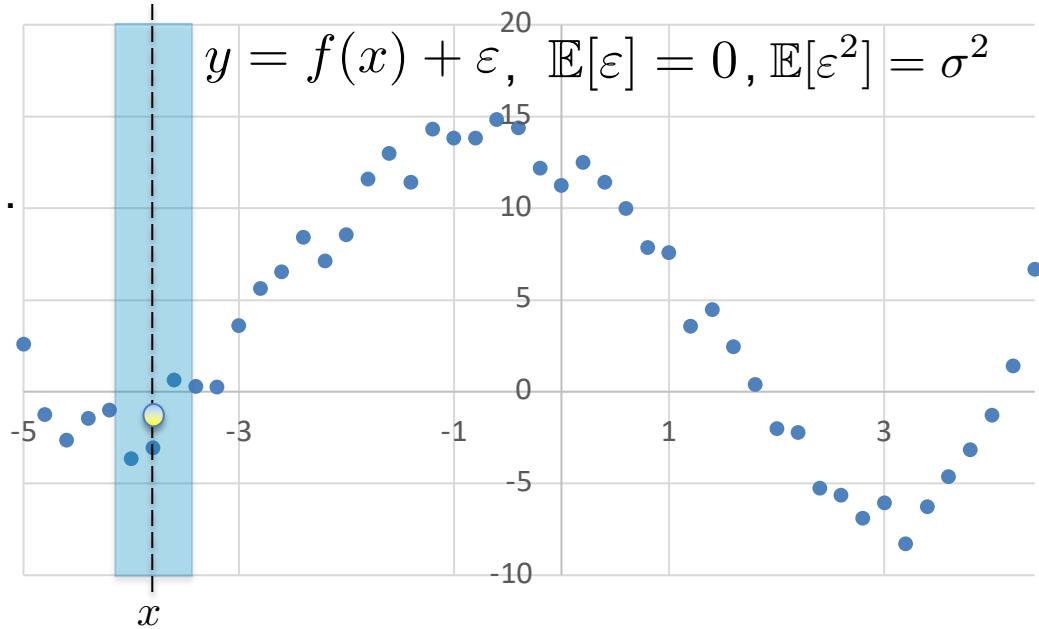
$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n.$$

ε_i i.i.d., $\mathbb{E}[\varepsilon_i] = 0$, $\mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$.

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in N_k(x)} y_i$$

Expected Prediction Error (EPE):

$$\mathbb{E} \left[(y - \hat{f}(x))^2 \right] = \mathbb{E} [(y - \mathbb{E}[y])^2] + (\mathbb{E}[y] - \mathbb{E}[\hat{f}(x)])^2 + \mathbb{E} \left[(\mathbb{E}[\hat{f}(x)] - \hat{f}(x))^2 \right]$$



Bias vs. Variance Trade-off

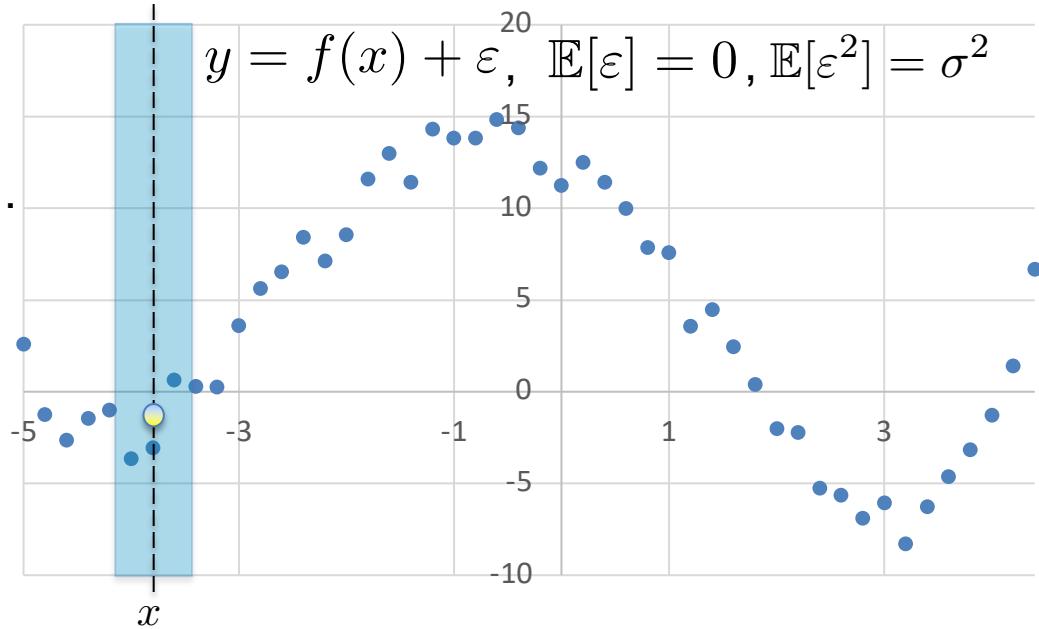
$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n.$$

$$\varepsilon_i \text{ i.i.d.}, \mathbb{E}[\varepsilon_i] = 0, \mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty.$$

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in N_k(x)} y_i$$

Expected Prediction Error (EPE):

$$\mathbb{E} \left[(y - \hat{f}(x))^2 \right] = \underbrace{\mathbb{E} [(y - \mathbb{E}[y])^2]}_{\text{inherent noise}} + (\mathbb{E}[y] - \mathbb{E}[\hat{f}(x)])^2 + \mathbb{E} \left[(\mathbb{E}[\hat{f}(x)] - \hat{f}(x))^2 \right]$$



Bias vs. Variance Trade-off

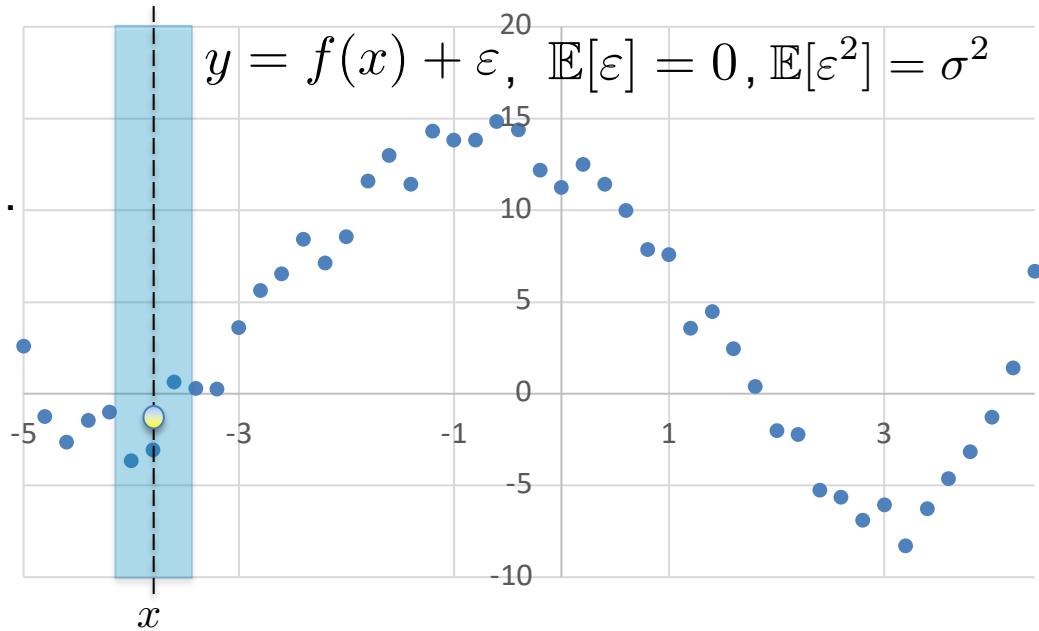
$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n.$$

$$\varepsilon_i \text{ i.i.d.}, \mathbb{E}[\varepsilon_i] = 0, \mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty.$$

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in N_k(x)} y_i$$

Expected Prediction Error (EPE):

$$\mathbb{E} \left[(y - \hat{f}(x))^2 \right] = \mathbb{E} [(y - \mathbb{E}[y])^2] + \underbrace{\left(\mathbb{E}[y] - \mathbb{E}[\hat{f}(x)] \right)^2}_{\text{estimator bias}} + \mathbb{E} \left[(\mathbb{E}[\hat{f}(x)] - \hat{f}(x))^2 \right]$$



Bias vs. Variance Trade-off

$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n.$$

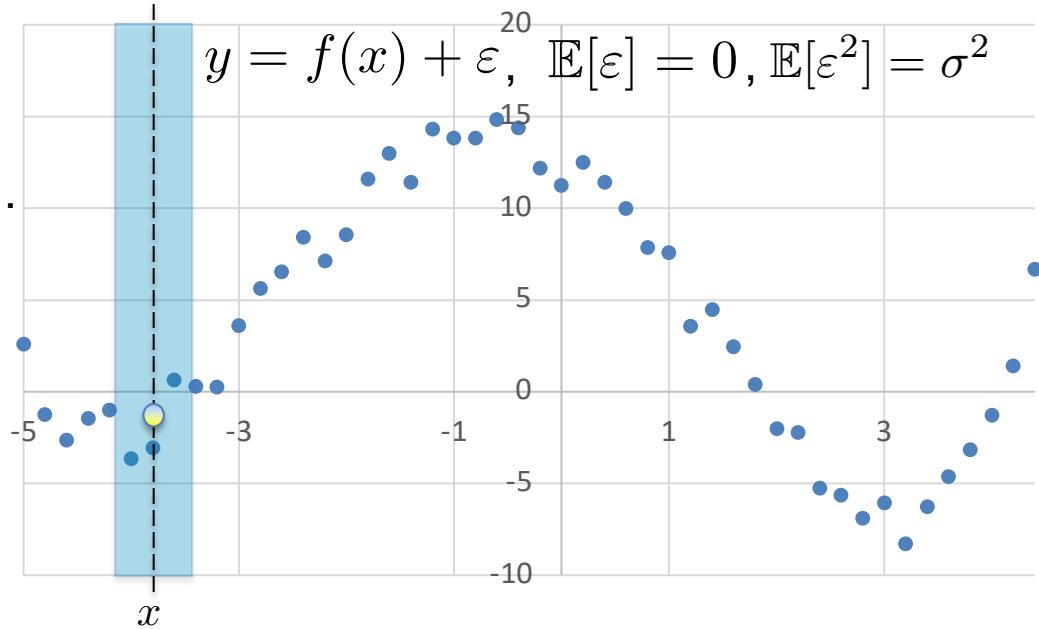
ε_i i.i.d., $\mathbb{E}[\varepsilon_i] = 0$, $\mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$.

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in N_k(x)} y_i$$

Expected Prediction Error (EPE):

$$\mathbb{E} \left[(y - \hat{f}(x))^2 \right] = \mathbb{E} [(y - \mathbb{E}[y])^2] + (\mathbb{E}[y] - \mathbb{E}[\hat{f}(x)])^2 + \mathbb{E} \left[(\mathbb{E}[\hat{f}(x)] - \hat{f}(x))^2 \right]$$

estimator variance



Bias vs. Variance Trade-off

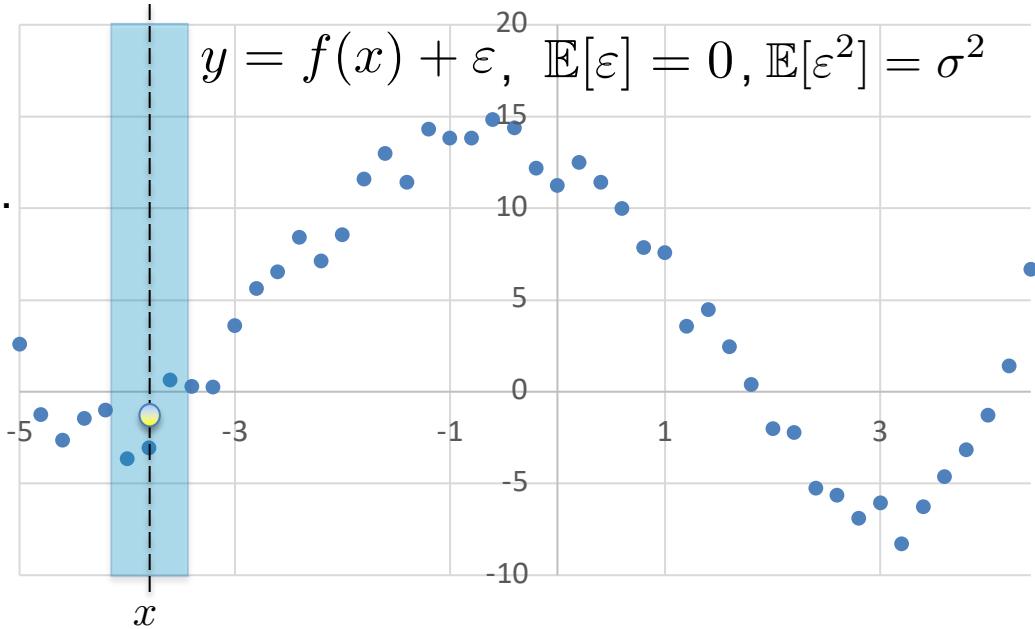
$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n.$$

$$\varepsilon_i \text{ i.i.d.}, \mathbb{E}[\varepsilon_i] = 0, \mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty.$$

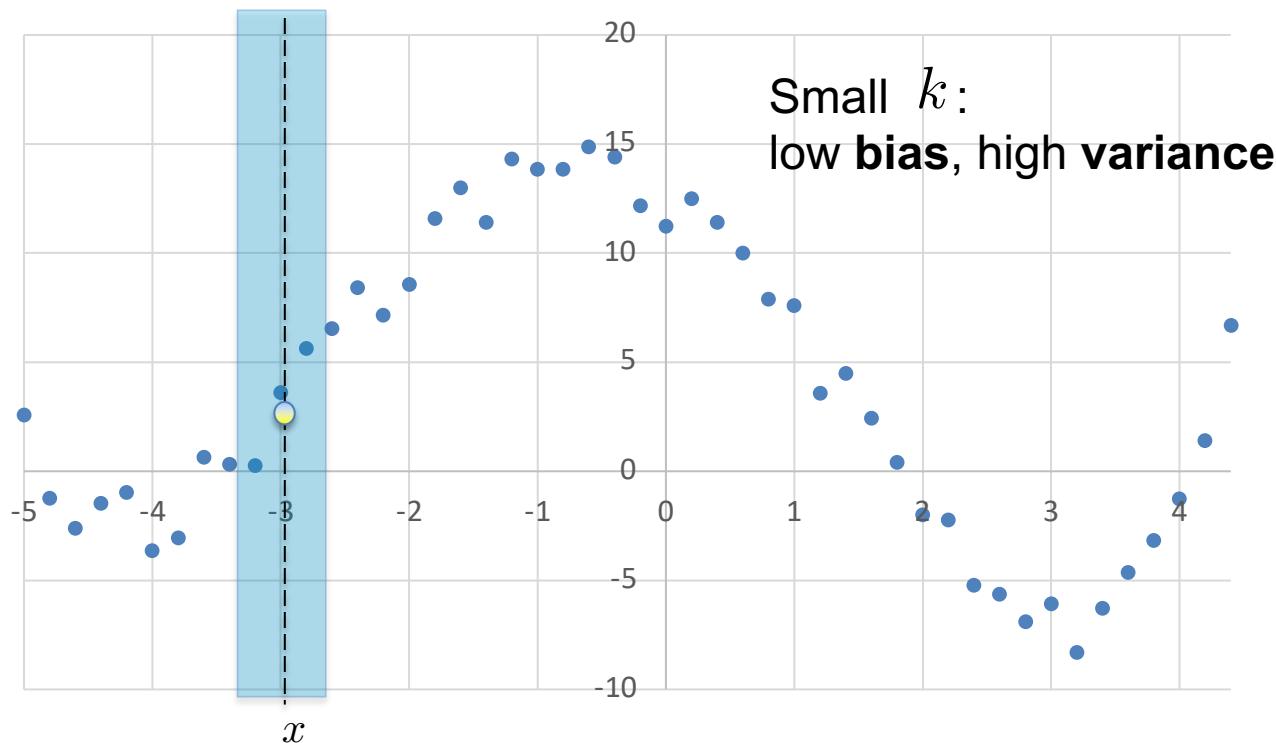
$$\hat{f}(x) = \frac{1}{k} \sum_{i \in N_k(x)} y_i$$

Expected Prediction Error (EPE):

$$\begin{aligned} \mathbb{E} \left[(y - \hat{f}(x))^2 \right] &= \mathbb{E} [(y - \mathbb{E}[y])^2] + (\mathbb{E}[y] - \mathbb{E}[\hat{f}(x)])^2 + \mathbb{E} \left[(\mathbb{E}[\hat{f}(x)] - \hat{f}(x))^2 \right] \\ &= \sigma^2 + \left(f(x) - \frac{1}{k} \sum_{i \in N_k(x)} f(x_i) \right)^2 + \frac{\sigma^2}{k} \end{aligned}$$

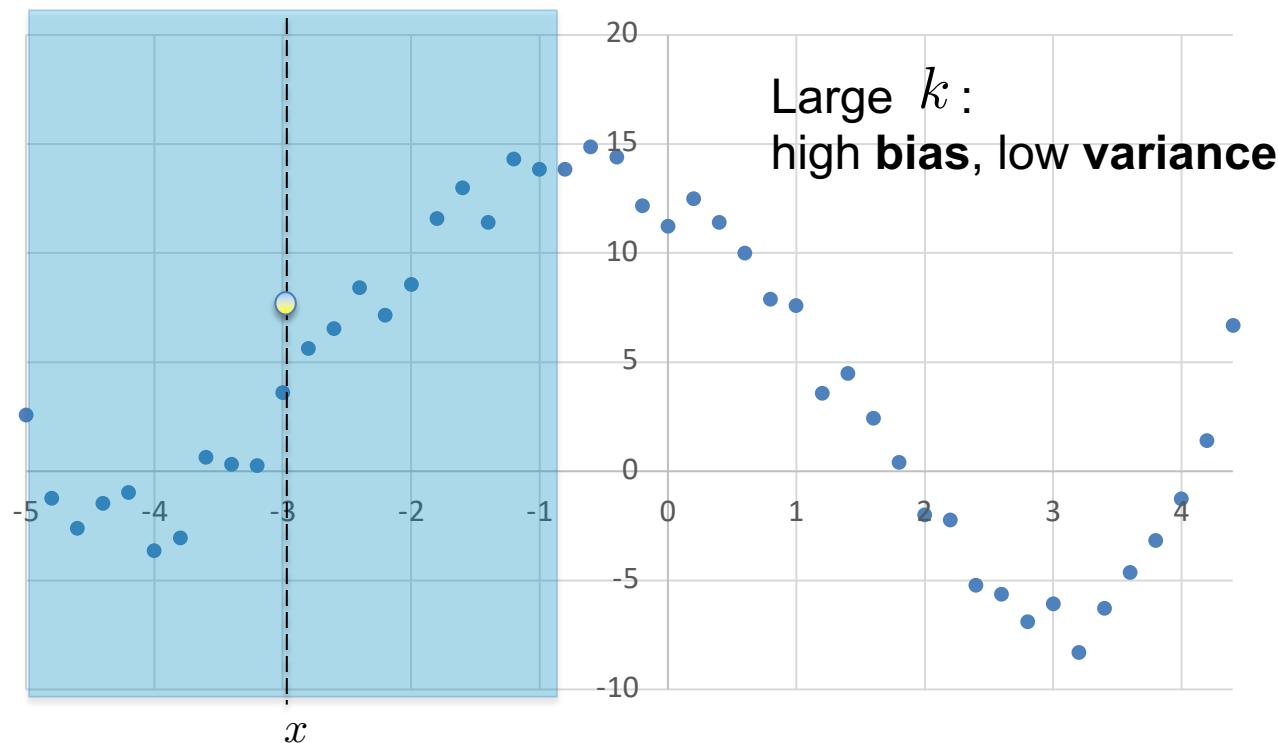


Bias vs. Variance Trade-off



$$\text{EPE: } \mathbb{E} \left[(y - \hat{f}(x))^2 \right] = \sigma^2 + \underbrace{\left(f(x) - \frac{1}{k} \sum_{i \in N_k(x)} f(x_i) \right)^2}_{\text{estimator bias}} + \underbrace{\frac{\sigma^2}{k}}_{\text{estimator variance}}$$

Bias vs. Variance Trade-off



$$\text{EPE: } \mathbb{E} \left[(y - \hat{f}(x))^2 \right] = \sigma^2 + \underbrace{\left(f(x) - \frac{1}{k} \sum_{i \in N_k(x)} f(x_i) \right)^2}_{\text{estimator bias}} + \underbrace{\frac{\sigma^2}{k}}_{\text{estimator variance}}$$

Bias vs. Variance Tradeoff

