

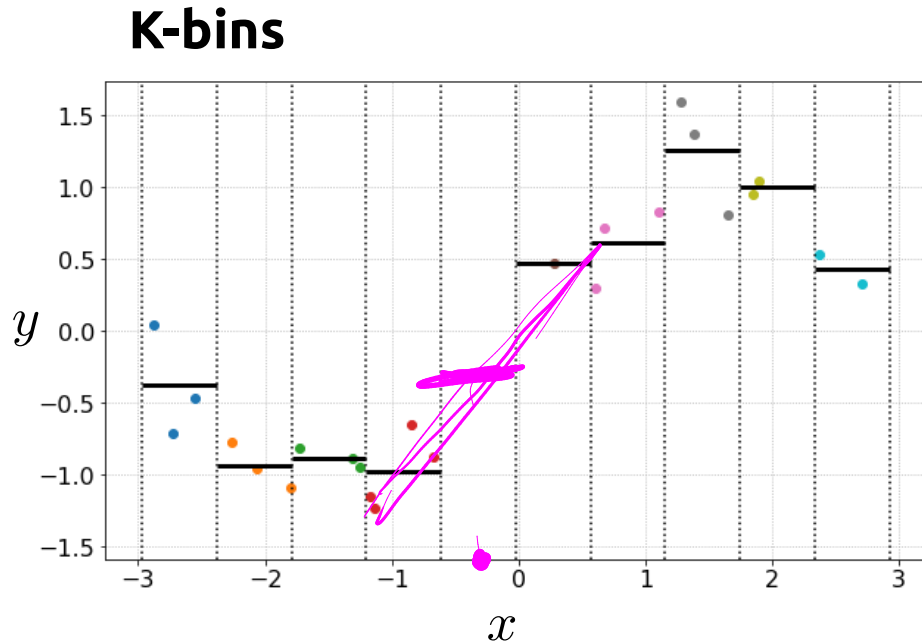


Statistics and Data Science for Engineers E178 / ME276DS

Linear regression Part 1

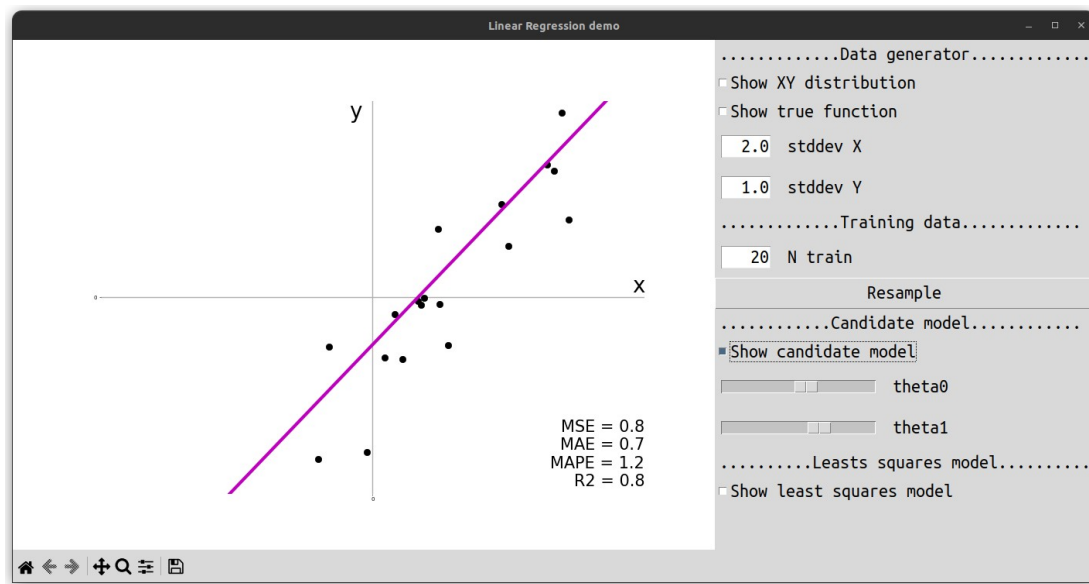
$$\mathcal{D} = \{(x_i, y_i)\}_{24}$$

x	y
2.704286	0.323832
-2.063888	-0.954098
0.606690	0.300653
-2.876493	0.043129
-1.725965	-0.816896
-1.174547	-1.156533
0.671117	0.720568
-1.247132	-0.946372
-1.801957	-1.087639
1.850384	0.952331
1.105398	0.829523
-2.267771	-0.778359
0.280262	0.472226
1.650797	0.808036
2.368964	0.530195
-2.728636	-0.707794
-0.667936	-0.880451
-1.314393	-0.891488
-2.552696	-0.468816
1.892769	1.047608
-0.849206	-0.653065
-1.134106	-1.231050
1.377637	1.367601
1.279469	1.590128

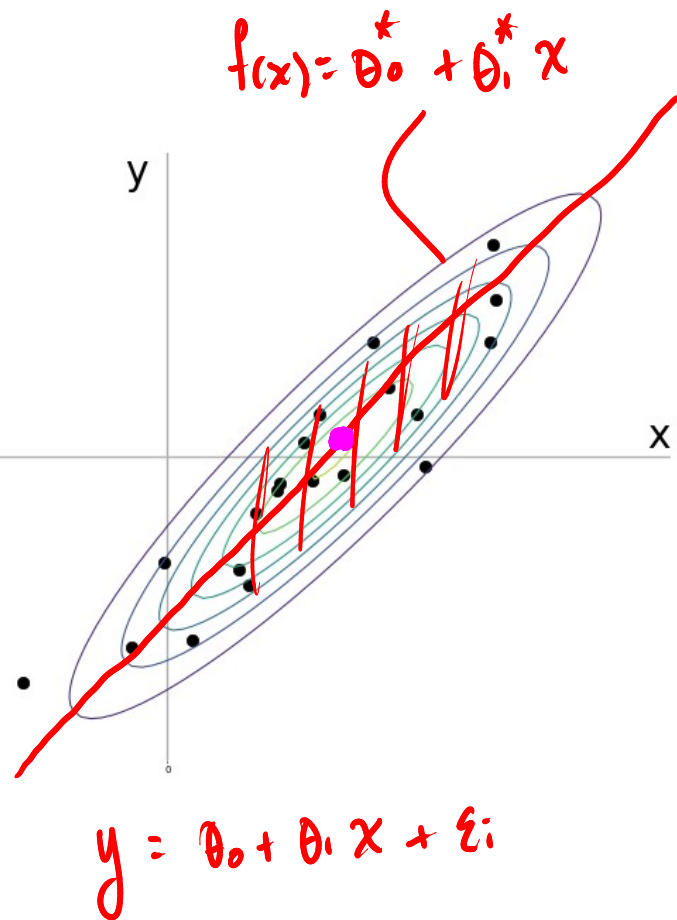
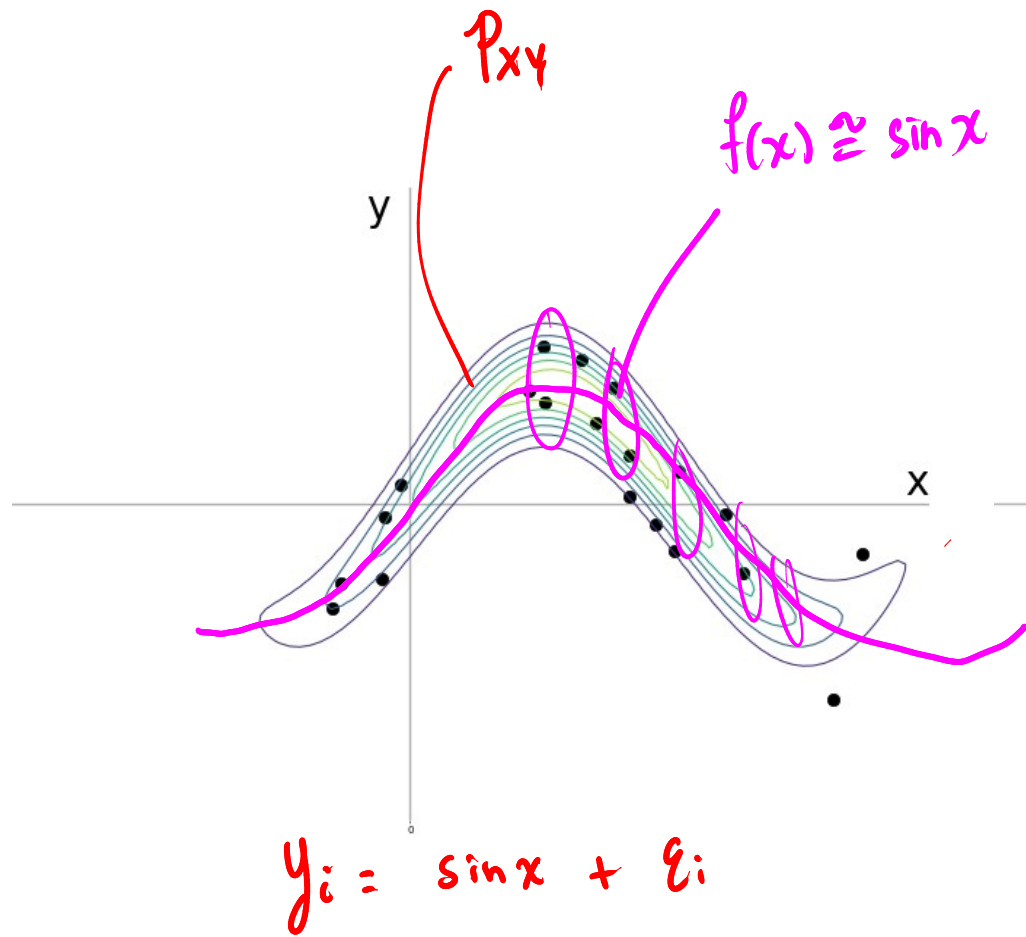


Interpolations.

Simple linear regression demo



$$y = \theta_0 + x \theta_1$$



$x_i \in \mathbb{R}, y_i \in \mathbb{R}.$

$\mathcal{D} = \{(x_i, y_i)\}_N$

Sample (point) estimates

$x \in \mathbb{R}^D \quad D=1.$

$$\begin{cases} \hat{\mu}_X = \frac{1}{N} \sum x_i & \dots E[X] \\ \hat{\mu}_Y = \frac{1}{N} \sum y_i & \dots E[Y]. \end{cases}$$

$$\hat{\sigma}_X^2 = \frac{1}{N-1} \sum (x_i - \hat{\mu}_X)^2 \dots \text{Var}[X]$$

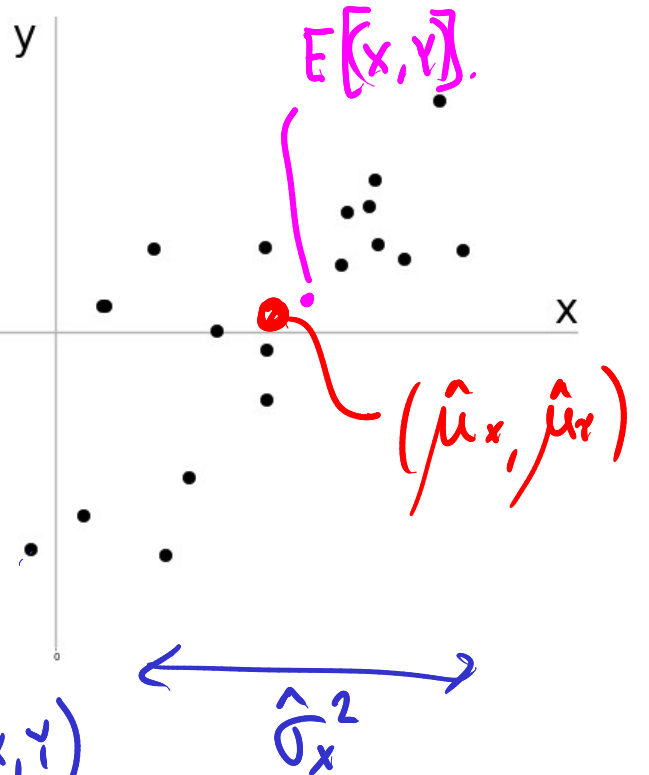
$$\hat{\sigma}_Y^2 = \frac{1}{N-1} \sum (y_i - \hat{\mu}_Y)^2 \dots \text{Var}[Y]$$

$$\hat{\sigma}_{XY} = \frac{1}{N-1} \sum (x_i - \hat{\mu}_X)(y_i - \hat{\mu}_Y) \dots \text{Cov}(X, Y)$$

$$r = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X \hat{\sigma}_Y} \in [-1, 1] \dots \text{sample correlation coefficient} \dots \rho_{XY}$$

All sums are over $i=1, \dots, N.$

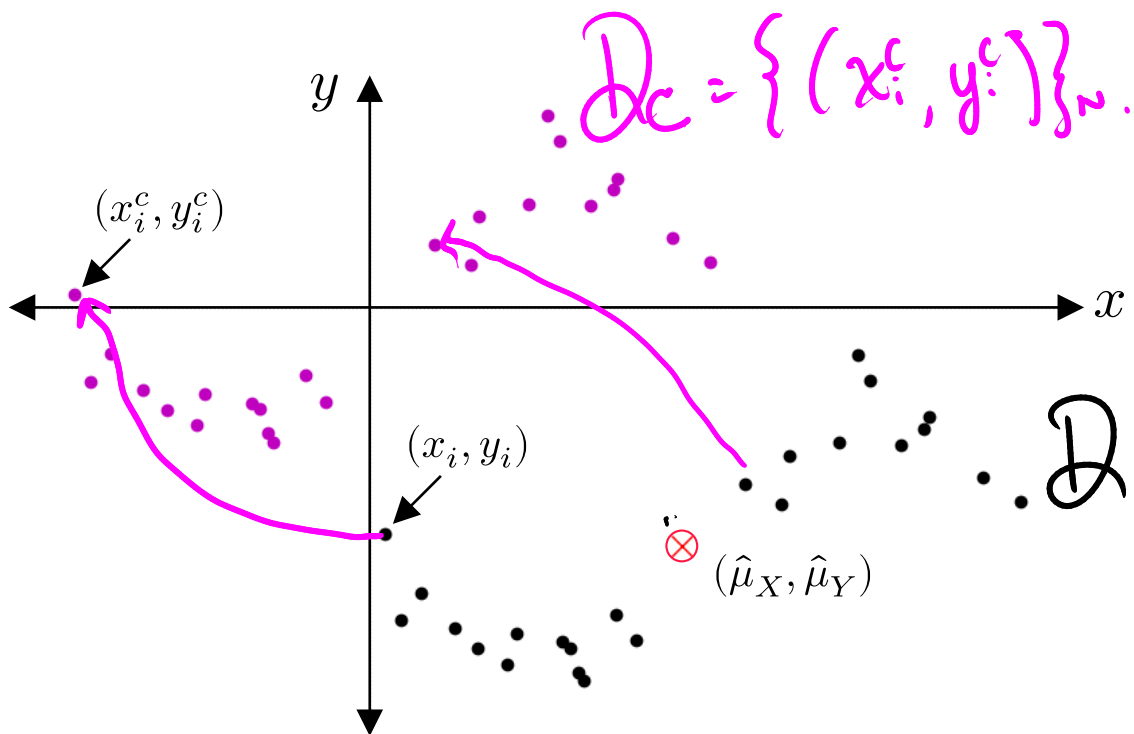
$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\text{Std}(X) \text{Std}(Y)}$$



Centering transformation

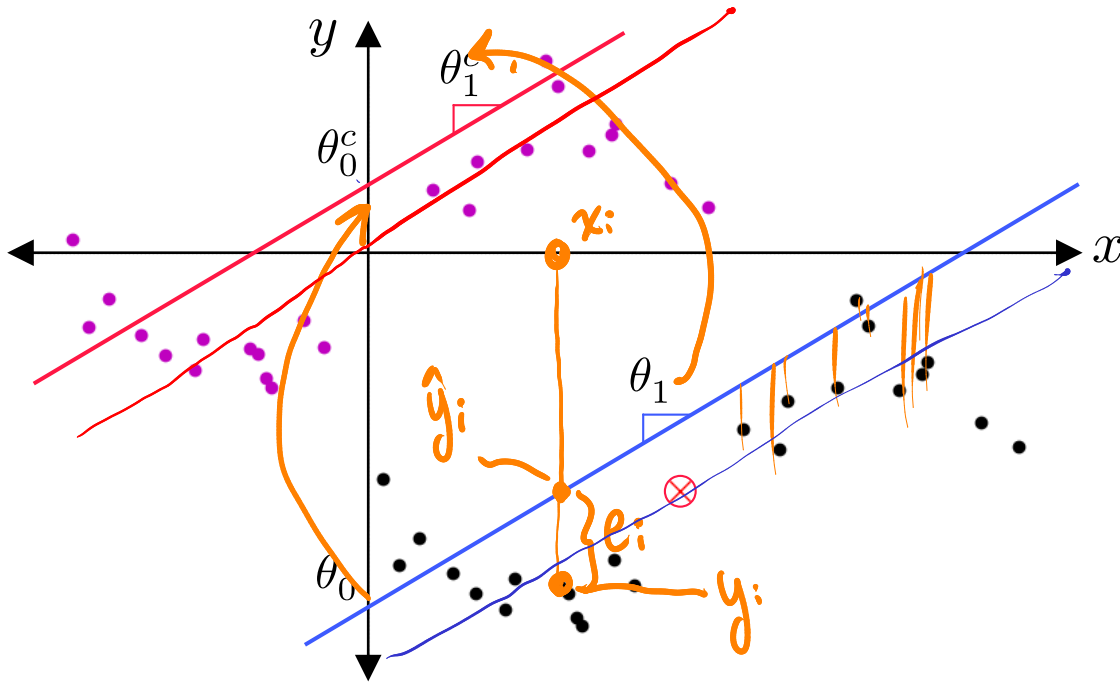
$$x_i^c = x_i - \hat{\mu}_X$$

$$y_i^c = y_i - \hat{\mu}_Y$$



Linear model

$$\hat{y}_i = h(x_i; \theta_0, \theta_1) = \theta_0 + x_i \theta_1 \quad i = 1 \dots N$$



Centered linear model

$$\hat{y}_i^c = \theta_0^c + x_i^c \theta_1^c$$

$$\begin{cases} \theta_1^c = \theta_1 \\ \theta_0^c = \theta_0 - \hat{\mu}_y + \hat{\mu}_x \theta_1 \end{cases}$$

optimal parameters.

Optimization problem

$$(\hat{\theta}_0, \hat{\theta}_1) = \operatorname{argmin}_{(\theta_0, \theta_1) \in \mathbb{R}^2}$$

$$= \operatorname{argmin}_{(\theta_0, \theta_1) \in \mathbb{R}^2}$$

$$= \operatorname{argmin}_{(\theta_0, \theta_1) \in \mathbb{R}^2}$$

$$\sum_{i=1}^N L_2(y_i, h(x_i; \theta_0, \theta_1))$$

$$\sum_{i=1}^N (\hat{y}_i - y_i)^2$$

$$\sum_{i=1}^N (\theta_0 + x_i \theta_1 - y_i)^2$$

squared

"least squares"

Centered problem:

$$(\hat{\theta}_0^c, \hat{\theta}_1^c) = \operatorname{argmin}_{(\theta_0, \theta_1) \in \mathbb{R}^2}$$

$$\sum_{i=1}^N (\theta_0 + x_i^c \theta_1 - y_i^c)^2$$

Stationary \Leftrightarrow global optimum.

Centered problem: Stationary points

$$J(\theta_0, \theta_1) = \sum_{i=1}^N (\theta_0 + x_i^c \theta_1 - y_i^c)^2$$

$$\frac{1}{N} \sum x_i^c = \hat{\mu}_x^c = 0,$$
$$\frac{1}{N} \sum y_i^c = \hat{\mu}_y^c = 0.$$

- $$\begin{aligned} \frac{\partial J}{\partial \theta_0} &= 2 \sum (\theta_0 + x_i^c \theta_1 - y_i^c) \\ &= 2 \left(N\theta_0 + \left(\sum x_i^c \right) \theta_1 - \sum y_i^c \right) \\ &= 2N\theta_0 = 0 \end{aligned}$$

$$\hat{\theta}_0^c = 0$$

⇒ L.S. regression
line goes
through
C.O.M. $(\hat{\mu}_x, \hat{\mu}_y)$.

- $$\begin{aligned} \frac{\partial J}{\partial \theta_1} &= 2 \sum (\theta_0 + x_i^c \theta_1 - y_i^c) x_i^c \\ &= 2 \left(\theta_0 \sum x_i^c + \theta_1 \sum (x_i^c)^2 - \sum x_i^c y_i^c \right) \end{aligned}$$

$$(N-1) \hat{\sigma}_{xy}$$

$$\sum (x_i^c)^2 = \sum (x_i - \mu_x)^2 = (N-1) \hat{\sigma}_x^2$$

$$\Rightarrow \theta_1 \hat{\sigma}_x^2 - \hat{\sigma}_{xy} = 0$$



$$\hat{\theta}_1^c$$

$$= \left[\frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x^2} = \hat{\theta}_1 \right]$$

Simple linear regression: Solution

$$\hat{\theta}_1 = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X^2}$$

$$\hat{\theta}_0 = \hat{\mu}_Y - \hat{\mu}_X \hat{\theta}_1$$

$$r = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X \hat{\sigma}_Y} = \hat{\theta}_1 \cdot \frac{\hat{\sigma}_X}{\hat{\sigma}_Y} \longrightarrow \boxed{\hat{\theta}_1 = r \frac{\hat{\sigma}_Y}{\hat{\sigma}_X}}$$

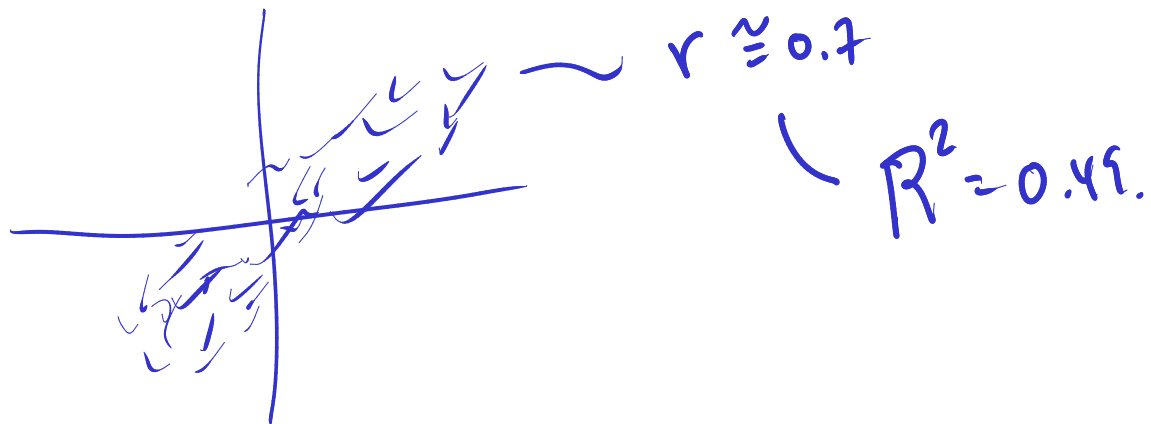
Simple linear regression: Performance

coefficient
of determination...

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \hat{\mu}_Y)^2}$$

$$= \vdots$$

$$= r^2 \quad \dots \text{square of the correlation.}$$



Statistical properties of simple linear regression

x 's (inputs) are fixed.

Assumed data

generating process:

$$y_i = \theta_0^* + x_i \theta_1^* + \varepsilon_i$$

linear function

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

Gaussian noise.

$$\left\{ \begin{array}{l} \hat{\theta}_0 = \bar{y} - \bar{x} \hat{\theta}_1 \\ \hat{\theta}_1 = r \frac{\hat{\sigma}_y}{\hat{\sigma}_x} \\ \hat{y}_i = \hat{\theta}_0 + \hat{\theta}_1 x_i \end{array} \right.$$

point estimators
of $\theta_0^*, \theta_1^*, y_i^*$

unbiased?
variance?

$\hat{\theta}_1$ as an estimate of θ_1^*

least
squares.

$$\hat{\theta}_1 = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X^2}$$

$$= \frac{1}{\hat{\sigma}_X^2} \frac{1}{N-1} \sum x_i^c y_i^c$$

$$= \frac{1}{(N-1)\hat{\sigma}_X^2} \sum x_i^c (\theta_1^* x_i^c + \varepsilon_i - \bar{\varepsilon})$$

$$= \frac{1}{(N-1)\hat{\sigma}_X^2} \left(\theta_1^* \sum (x_i^c)^2 + \sum x_i^c \varepsilon_i - \bar{\varepsilon} \sum x_i^c \right)$$

$$= \frac{\theta_1^*}{\hat{\sigma}_X^2} \frac{\sum (x_i^c)^2}{(N-1)} + \frac{\sum x_i^c \varepsilon_i}{(N-1)\hat{\sigma}_X^2}$$

$$= \theta_1^* + \frac{\sum x_i^c \varepsilon_i}{(N-1)\hat{\sigma}_X^2}$$

everything is deterministic
except the ε_i .

Define: $\bar{\varepsilon} = \frac{1}{N} \sum_{j=1}^N \varepsilon_j$

Use: $y_i^c = x_i^c \theta_1^* + \varepsilon_i - \bar{\varepsilon}$
(proved in the reader)

$$\hat{\theta}_1 = \hat{\theta}_1$$

$$\varepsilon_i = \varepsilon_i$$

$\hat{\theta}_1$ as an estimate of θ_1^*

$$\hat{\theta}_1 = \theta_1^* + \frac{\sum x_i^c \varepsilon_i}{(N-1)\hat{\sigma}_X^2}$$

substitute
in R.V.s.

\Rightarrow

$$\hat{\Theta}_1 = \theta_1^* + \frac{\sum x_i \varepsilon_i}{(N-1)\hat{\sigma}_X^2}$$

$$E[\hat{\Theta}_1] = E\left[\theta_1^* + \frac{\sum x_i^c \varepsilon_i}{(N-1)\hat{\sigma}_X^2}\right]$$

$$= \theta_1^* + \frac{\sum x_i^c E[\varepsilon_i]}{(N-1)\hat{\sigma}_X^2}$$

$$= \theta_1^*$$

... $\hat{\Theta}_1$ is an unbiased estimate of θ_1^*

$\hat{\theta}_1$ as an estimate of θ_1^*

$$Var[\hat{\Theta}_1] = Var \left[\theta_1^* + \frac{\sum x_i^c \mathcal{E}_i}{(N-1)\hat{\sigma}_X^2} \right]$$

$$= \sum \left(\frac{(x_i^c)^2}{(N-1)^2 \hat{\sigma}_X^4} Var[\mathcal{E}_i] \right)$$

$$= \frac{\sigma^2}{(N-1)^2 \hat{\sigma}_X^4} \sum (x_i^c)^2$$

$$= \frac{\sigma^2}{(N-1)\hat{\sigma}_X^2}$$

vertical
variance

horizontal
spread.

$\hat{\theta}_0$ as an estimate of θ_0^*

Express in terms
of well described
quantities:

$$\hat{\theta}_0 = \theta_0^* + \hat{\mu}_X \theta_1^* + \bar{\varepsilon} - \hat{\mu}_X \hat{\theta}_1$$

(Handwritten red checkmarks are above each term in the equation.)

$$\bar{\varepsilon} = \frac{1}{N} \sum \varepsilon_i$$
$$E[\bar{\varepsilon}] = 0 \quad \text{Var}[\bar{\varepsilon}] = \frac{\sigma^2}{N}$$

(proved in the reader)

Substitute RVs:

$$\hat{\Theta}_0 = \theta_0^* + \hat{\mu}_X \theta_1^* + \bar{\varepsilon} - \hat{\mu}_X \hat{\Theta}_1$$

(Handwritten red arrows point from the corresponding terms in the equation above to this one.)

Statistical analysis:

$$E[\hat{\Theta}_0] = E[\theta_0^* + \hat{\mu}_X \theta_1^* + \bar{\varepsilon} - \hat{\mu}_X \hat{\Theta}_1]$$
$$= \theta_0^* \quad \dots \text{unbiased!}$$

(proved in the reader)

$$\text{Var}[\hat{\Theta}_0] = \frac{\sigma^2}{N} + \frac{\sigma^2 \hat{\mu}_X^2}{(N-1) \hat{\sigma}_X^2}$$

(Handwritten red arrow points to the denominator term $\hat{\sigma}_X^2$.)

(proved in the reader)

\hat{y}_i as an estimate of y_i^*

Express in terms
of well described
quantities:

$$\hat{y}_i = \theta_0^* + \hat{\mu}_X \theta_1^* + \bar{\varepsilon} + x_i^c \hat{\theta}_1 \quad (\text{proved in the reader})$$

rv.

Substitute RVs:

$$\hat{Y}_i = \theta_0^* + \hat{\mu}_X \theta_1^* + \bar{\mathcal{E}} + x_i^c \hat{\Theta}_1$$

↓

Statistical analysis:

$$E[\hat{Y}_i] = \theta_0^* + x_i \theta_1^* \quad \dots \text{also unbiased!} \quad (\text{proved in the reader})$$

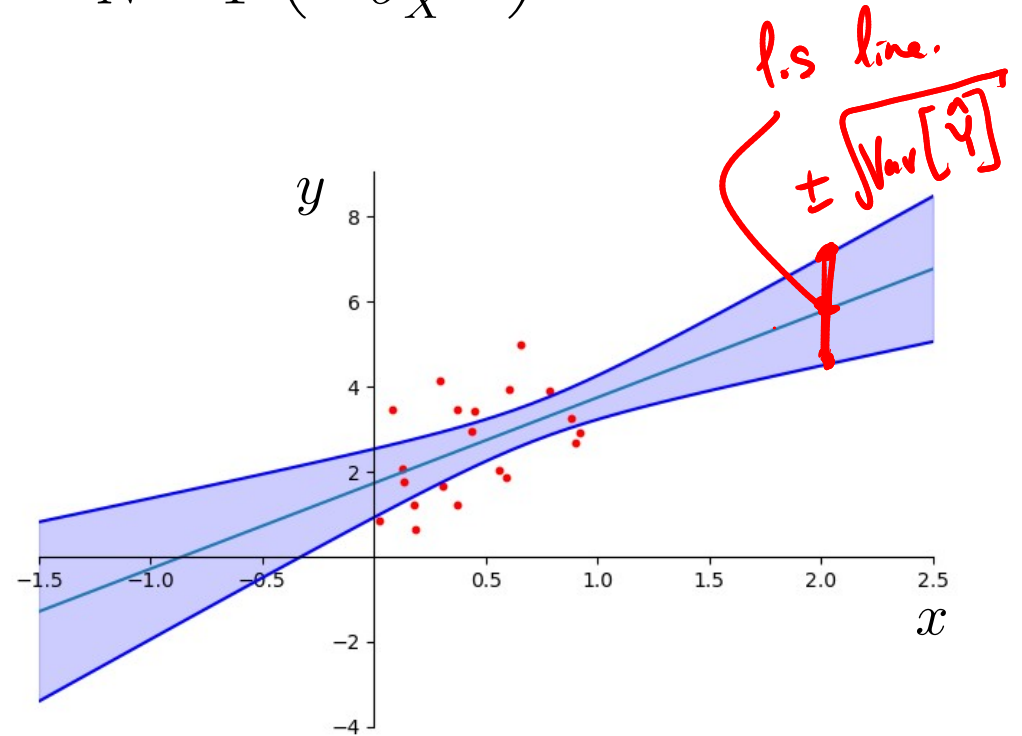
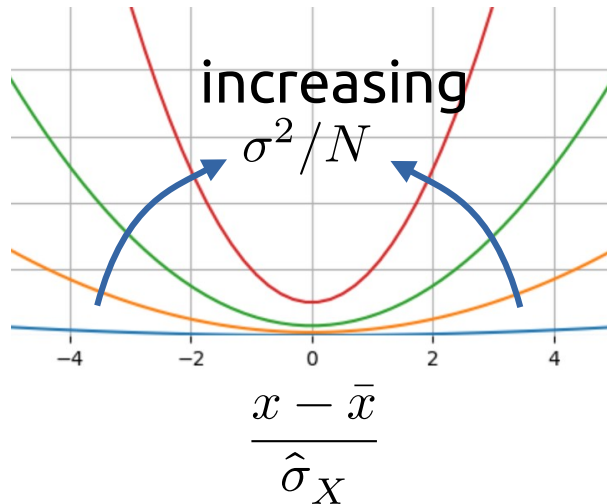
$$Var[\hat{Y}_i] = \frac{\sigma^2}{N} + \frac{\sigma^2 (x_i - \hat{\mu}_X)^2}{\sum_{j=1}^N (x_j - \hat{\mu}_X)^2} \quad (\text{proved in the reader})$$

$(n-1) \hat{\sigma}_x^2$

rewrite

Prediction uncertainty in simple linear regression

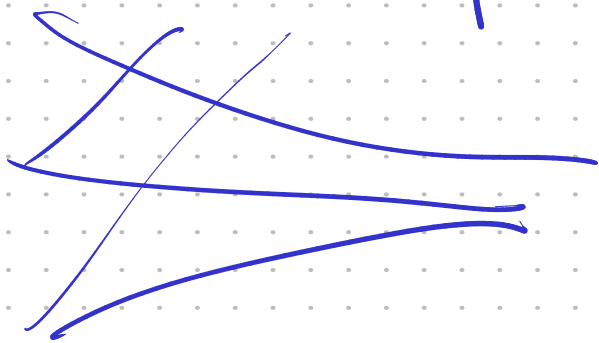
$$\text{Stdev}[\hat{Y}] = \sqrt{\text{Var}[\hat{Y}] = \frac{\sigma^2}{N} + \frac{\sigma^2}{N-1} \left(\frac{x - \hat{\mu}_X}{\hat{\sigma}_X} \right)^2}$$



$\hat{\theta}_0$: know expected value and variance

$\hat{\theta}_1$

\hat{y}_i



Gaussian!

know the distributions
of these estimators!

Confidence intervals and hypothesis tests

- Find confidence interval for the parameters
- Find confidence interval for the output
- Hypothesis test: $H_0 : \theta_d = 0$
 $H_1 : \theta_d \neq 0$
- Unknown σ^2 :

If unknown:

$$\hat{\sigma}^2 = \frac{1}{N - \cancel{D}} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

and use t-test.

~~1~~ 2

90% confidence

$$\hat{\theta}_0 \in [0.2, 0.3]$$

$$\hat{\theta}_1 \in [0.7, 2.2]$$

$$\hat{y}_i = \text{---} \pm \text{---}$$

$$\cancel{d=1} \dots \cancel{D} \quad D=1$$

"y does not depend on x"

"y does depend on x"

$$= \frac{1}{N-2} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$D > 1$

Multiple linear regression ($D > 1$)

$$\mathcal{D} = \{(x_i, y_i)\}_N = \{(x_i^1, \dots, x_i^D, y_i)\}_N$$

$$x_i = [x_i^1 \quad \dots \quad x_i^D] \in \mathbb{R}^{1 \times D}$$

$$\mathbf{X} = \begin{bmatrix} x_1^1 & \dots & x_1^D \\ \vdots & & \vdots \\ x_N^1 & \dots & x_N^D \end{bmatrix} \in \mathbb{R}^{N \times D}$$

$$\mathbf{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^{N \times 1}$$

	x^1	x^2	x^3	...	x^D	y
0	0.247746	36.0	266.0	0.247746	88.019654	57.000823
1	0.179340	37.0	365.0	0.179340	27.211935	22.471227
2	0.956807	21.0	151.0	0.956807	97.456012	56.357366
3	0.869653	43.0	437.0	0.869653	43.203221	34.59475
4	0.825345	47.0	160.0	0.825345	98.933930	64.426163
5	0.331114	321.0	371.0	0.331114	8.257917	14.218720
6	0.765523	17.0	364.0	0.765523	96.696783	59.123769
7	0.956807	21.0	151.0	0.956807	97.456012	52.983470

$\hat{\mu}_X$

$\hat{\mu}_Y$

$$\hat{\mu}_X = \frac{1}{N} \mathbf{1}_N^T \mathbf{X}$$

$$\hat{\mu}_Y = \frac{1}{N} \mathbf{1}_N^T \mathbf{Y}$$

Data-centering transformation:

$$\begin{aligned}\mathbf{X}^c &= \mathbf{X} - \mathbf{1}_N \hat{\mu}_X & \mathbf{1}_N^T \mathbf{X}^c &= \mathbf{0}_D \\ \mathbf{Y}^c &= \mathbf{Y} - \mathbf{1}_N \hat{\mu}_Y & \mathbf{1}_N^T \mathbf{Y}^c &= 0\end{aligned}$$

Linear model:

$$\begin{aligned}\hat{y}_i &= \theta_0 + x_i^1 \theta_1 + \dots + x_i^D \theta_D \\ &= \theta_0 + x_i \underline{\theta}_1\end{aligned}\quad \underline{\theta}_1 = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta^D \end{bmatrix} \in \mathbb{R}^{D \times 1}$$

Matrix form:

$$\hat{\mathbf{Y}} = \mathbf{1}_N \theta_0 + \mathbf{X} \underline{\theta}_1$$

Centered:

$$\hat{\mathbf{Y}}^c = \mathbf{1}_N \theta_0^c + \mathbf{X}^c \underline{\theta}_1^c$$

Optimization problem

Original: $(\hat{\theta}_0, \hat{\underline{\theta}}_1) = \underset{(\theta_0, \underline{\theta}_1) \in \mathbb{R}^{D+1}}{\operatorname{argmin}} \sum_{i=1}^N (\theta_0 + x_i \underline{\theta}_1 - y_i)^2$

Matrix: $(\hat{\theta}_0, \hat{\underline{\theta}}_1) = \underset{(\theta_0, \underline{\theta}_1) \in \mathbb{R}^{D+1}}{\operatorname{argmin}} (\mathbf{1}_N \theta_0 + \mathbf{X} \underline{\theta}_1 - \mathbf{Y})^T (\mathbf{1}_N \theta_0 + \mathbf{X} \underline{\theta}_1 - \mathbf{Y})$

Centered: $(\hat{\theta}_0^c, \hat{\underline{\theta}}_1^c) = \underset{(\theta_0, \underline{\theta}_1) \in \mathbb{R}^{D+1}}{\operatorname{argmin}} (\mathbf{1}_N \theta_0 + \mathbf{X}^c \underline{\theta}_1 - \mathbf{Y}^c)^T (\mathbf{1}_N \theta_0 + \mathbf{X}^c \underline{\theta}_1 - \mathbf{Y}^c)$

Condensed: $(\hat{\theta}_0^c, \hat{\theta}_1^c) = \underset{\underline{\theta}^c \in \mathbb{R}^{D+1}}{\operatorname{argmin}} \quad (\mathbb{X} \underline{\theta}^c - \mathbf{Y}^c)^T (\mathbb{X} \underline{\theta}^c - \mathbf{Y}^c)$

where $\mathbb{X} = [\mathbf{1}_N \quad \mathbf{X}^c] \quad \underline{\theta}^c = \begin{bmatrix} \theta_0^c \\ \theta_1^c \end{bmatrix}$

Cost function:
$$\begin{aligned} J(\underline{\theta}^c) &= (\mathbb{X} \underline{\theta}^c - \mathbf{Y}^c)^T (\mathbb{X} \underline{\theta}^c - \mathbf{Y}^c) \\ &= ((\underline{\theta}^c)^T \mathbb{X}^T - (\mathbf{Y}^c)^T) (\mathbb{X} \underline{\theta}^c - \mathbf{Y}^c) \\ &= (\underline{\theta}^c)^T \mathbb{X}^T \mathbb{X} \underline{\theta}^c - 2(\mathbf{Y}^c)^T \mathbb{X} \underline{\theta}^c + (\mathbf{Y}^c)^T \mathbf{Y}^c \end{aligned}$$

Gradient: $\nabla J = 2\mathbb{X}^T \mathbb{X} \underline{\theta}^c - 2\mathbb{X}^T \mathbf{Y}^c$

Stationary points: $\hat{\underline{\theta}}^c = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{Y}^c$

Unpack:

$$\begin{aligned} \begin{bmatrix} \hat{\theta}_0^c \\ \hat{\theta}_1^c \end{bmatrix} &= \left(\begin{bmatrix} \mathbf{1}_N^T \\ (\mathbf{X}^c)^T \end{bmatrix} \begin{bmatrix} \mathbf{1}_N & \mathbf{X}^c \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{1}_N^T \\ (\mathbf{X}^c)^T \end{bmatrix} \mathbf{Y}^c \\ &= \begin{bmatrix} \vdots \\ 0 \\ ((\mathbf{X}^c)^T \mathbf{X}^c)^{-1} (\mathbf{X}^c)^T \mathbf{Y}^c \end{bmatrix} \end{aligned}$$

Un-center:

$$\begin{aligned} \hat{\theta}_0 &= \hat{\mu}_Y - \hat{\mu}_X \hat{\theta}_1 \\ \hat{\theta}_1 &= ((\mathbf{X}^c)^T \mathbf{X}^c)^{-1} (\mathbf{X}^c)^T \mathbf{Y}^c \end{aligned}$$

Statistical properties of multiple linear regression

Assumed data
generating process:

$$y_i = \theta_0^* + x_i \theta_1^* + \varepsilon_i$$
$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

Centered process:

$$y_i^c = (\theta_0^*)^c + x_i^c (\theta_1^*)^c + \varepsilon_i^c$$
$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$\hat{\underline{\theta}}_1$ as an estimate of $\underline{\theta}_1^*$

$$\hat{\underline{\theta}}_1 = \underline{\theta}_1^* + \Sigma_X^{-1}(\mathbf{X}^c)^T \underline{\varepsilon}^c$$

where $\Sigma_X^{-1} = ((\mathbf{X}^c)^T \mathbf{X}^c)^{-1}$

$$\hat{\underline{\Theta}}_1 = \underline{\theta}_1^* + \Sigma_X^{-1}(\mathbf{X}^c)^T \underline{\mathcal{E}}^c$$

(proved in the reader)

$$E [\hat{\underline{\Theta}}_1] = \underline{\theta}_1^*$$

$$Var [\hat{\underline{\Theta}}_1] = \sigma^2 \Sigma_X^{-1}$$

$\hat{\theta}_0$ as an estimate of θ_0^*

$$\hat{\theta}_0 = \theta_0^* + \hat{\mu}_X \theta_1^* + \bar{\varepsilon} - \hat{\mu}_X \hat{\theta}_1$$

(proved in the reader)

$$\hat{\Theta}_0 = \theta_0^* + \hat{\mu}_X \theta_1^* + \bar{\varepsilon} - \hat{\mu}_X \hat{\Theta}_1$$

$$E [\hat{\Theta}_0] = \theta_0^*$$

$$Var [\hat{\Theta}_0] = \frac{\sigma^2}{N} + \sigma^2 \hat{\mu}_X \Sigma_X^{-1} \hat{\mu}_X^T$$

\hat{y}_i as an estimate of y_i^*

Same as the scalar case:

$$\hat{Y}_i = \theta_0^* + \hat{\mu}_X \theta_1^* + \bar{\mathcal{E}} + x_i^c \hat{\Theta}_1$$

$$E[\hat{Y}_i] = \theta_0^* + \theta_1^* x_i$$

$$Var[\hat{Y}_i] = \frac{\sigma^2}{N} + \sigma^2 (x_i - \hat{\mu}_X) \Sigma_X^{-1} (x_i - \hat{\mu}_X)^T$$