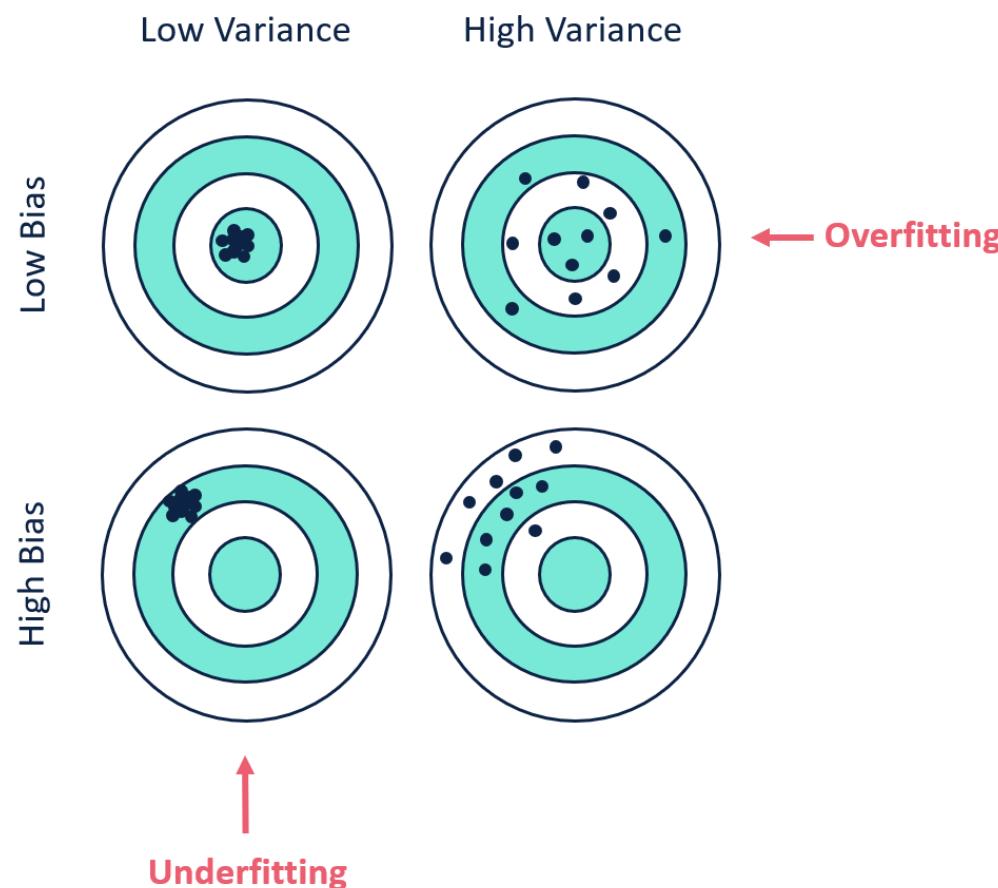
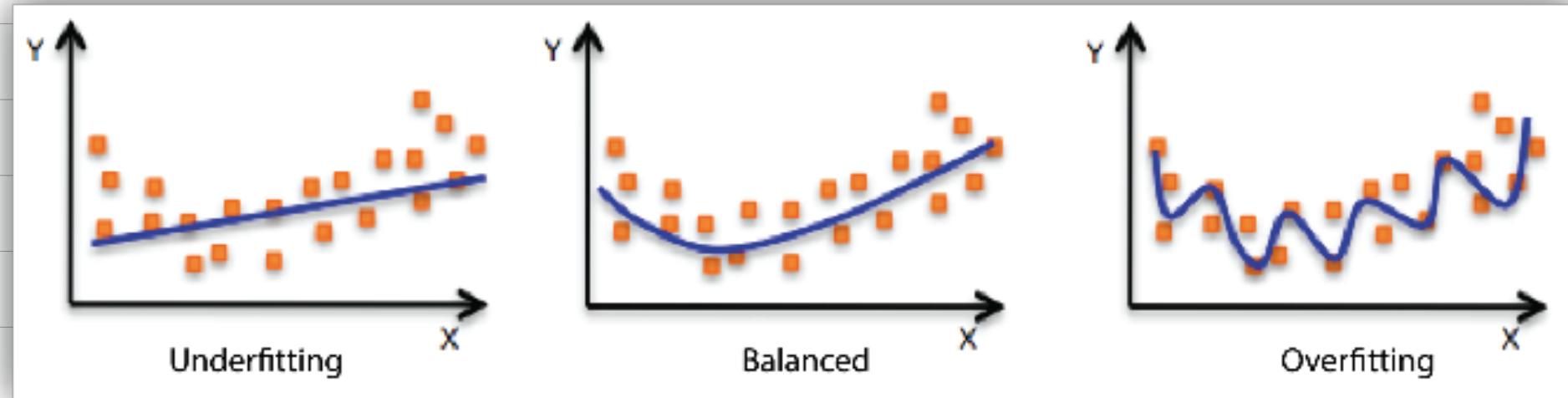


Session 6

Regularization





عدم دقت حاصل از تخمین یک مدل های واقعی **biase**

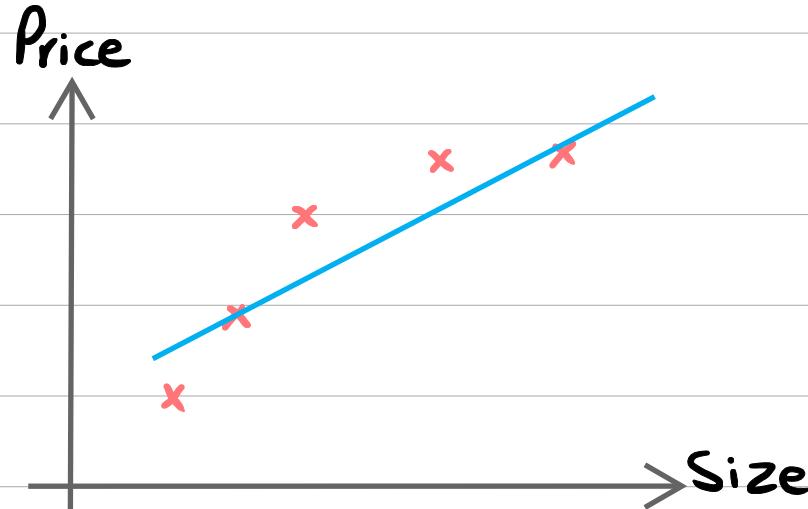
با چیزی زیاد با یک مدل ساده.

train او را با یک گروه داده $f_{w,b}(x)$ ار: **Variance**

متغیرات تخمین بزرگ، چند تغییر می کند.

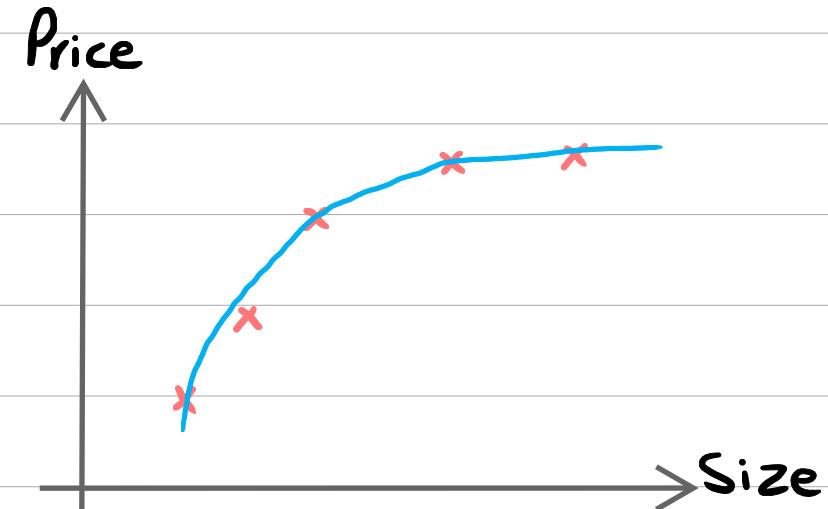
$$\text{Error} = \text{var}(f(x)) + \text{Bias}(f(x))^2 + \text{error}(\epsilon)$$

Regression



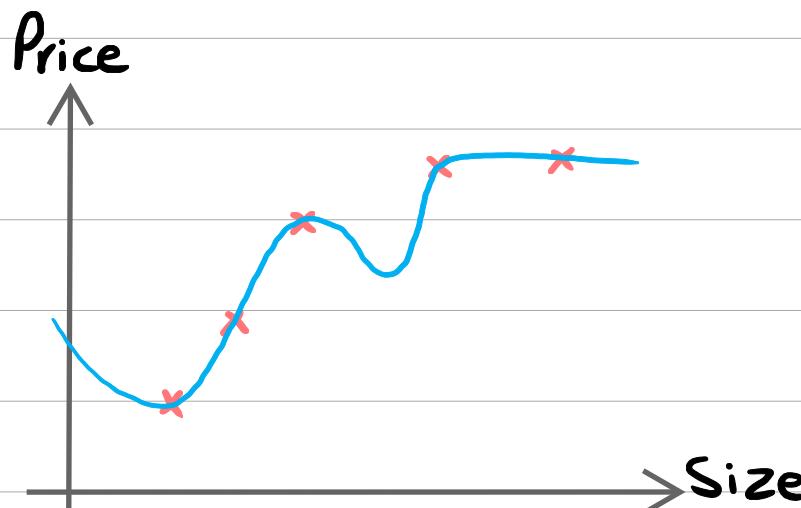
$$w_1x + b$$

Underfit \rightsquigarrow high bias \uparrow



$$w_1x + w_2x^2 + b$$

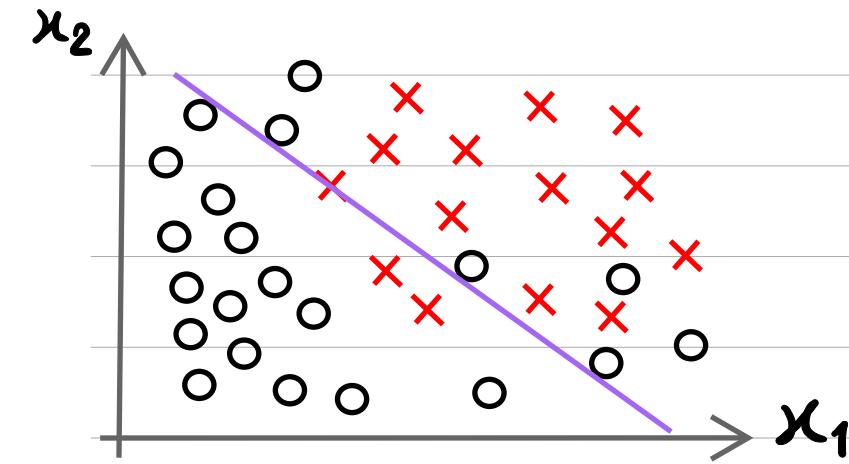
Just right \rightsquigarrow generalization



$$w_1x + w_2x^2 + w_3x^3 + w_4x^4 + b$$

Overfit \rightsquigarrow high variance \uparrow

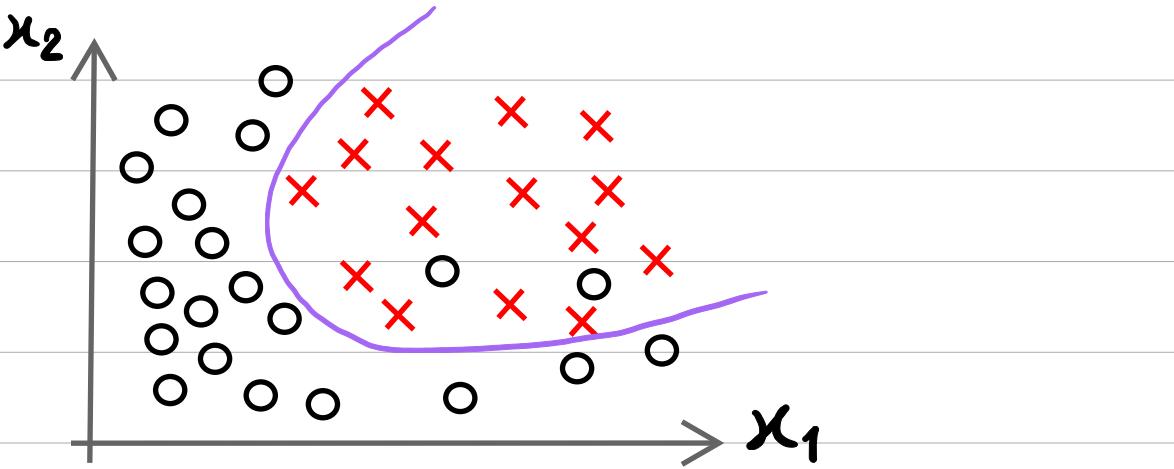
Classification



$$Z = w_1 x_1 + w_2 x_2 + b$$

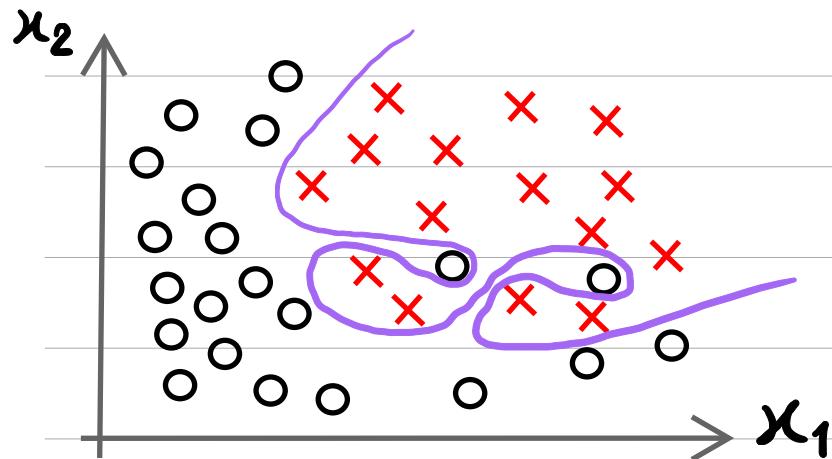
$$f_{w,b}(x) = g(z) \rightarrow g: \text{Sigmoid}$$

Underfit \rightsquigarrow high bias \uparrow



$$\begin{aligned} Z = & w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 \\ & + w_5 x_1 x_2 + b \end{aligned}$$

Just right \rightsquigarrow generalization

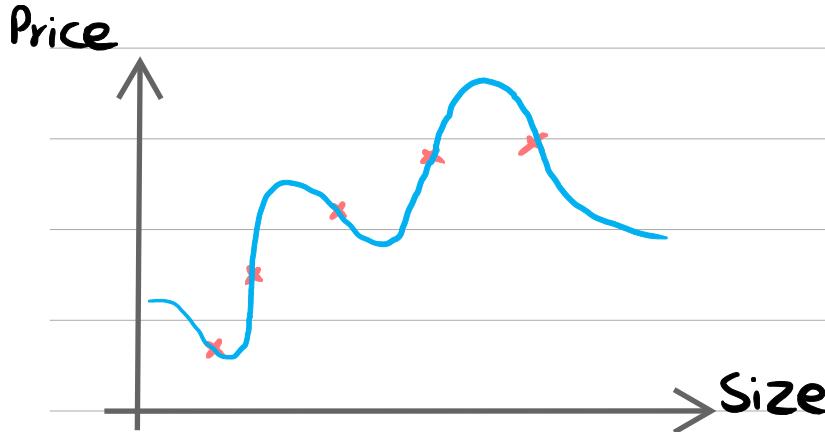


$$\begin{aligned} Z = & w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2 \\ & + w_5 x_1^2 x_2^3 + w_6 x_1^3 x_2 + \dots + b \end{aligned}$$

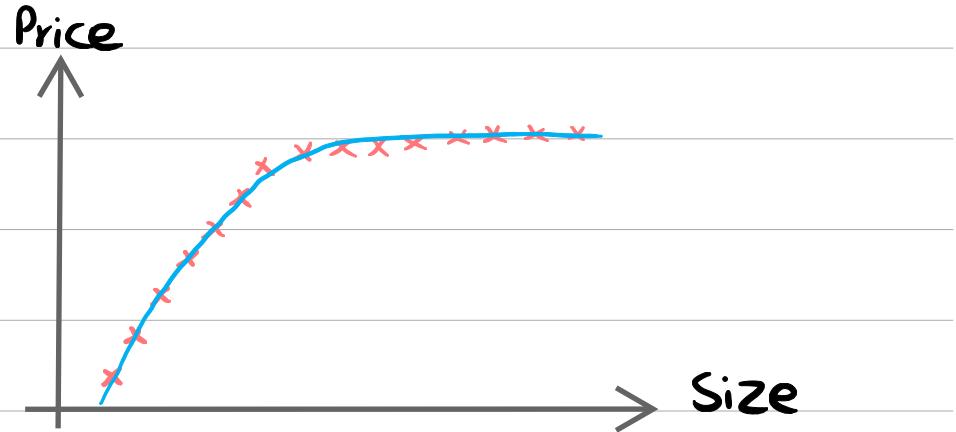
Overfit \rightsquigarrow high variance \uparrow

1. Collect more training examples

؟ برای رفع مشکل overfit کنیم؟



Overfit



Good fit

2. Select features to include/exclude

x_1	x_2	x_3	x_4	x_5	...	x_{100}	y
Size	bedroom	floors	age	avg income	...	distance coffee shop	Price

all features

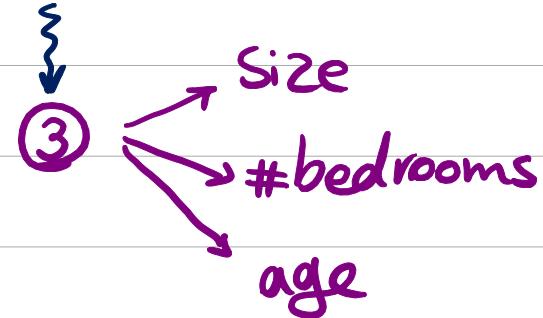
+

Insufficient sample

Overfit

کافی نبودن ممونها

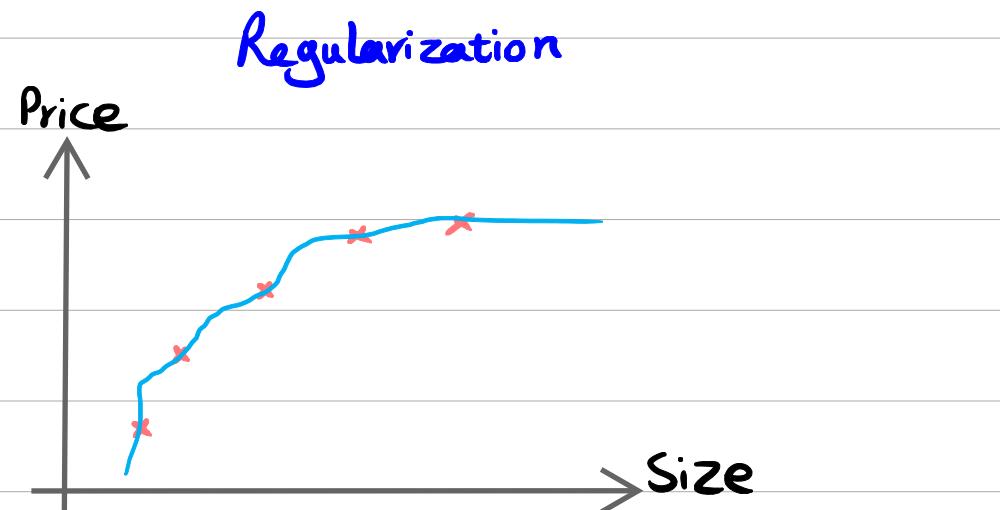
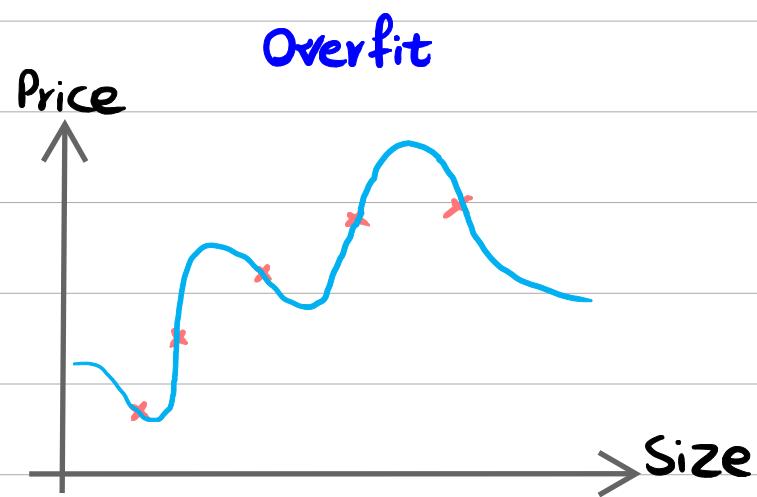
Selected features



مسکن است دیگر که های

نمی بذف شوند

3. Reduce the size of parameters W_j



$$f(x) = 28x - 385x^2 + 39x^3$$

$$-174x^4 + 100$$

Large values for W_j

$$f(x) = 13x - 0.23x^2 + 0.000014x^3$$

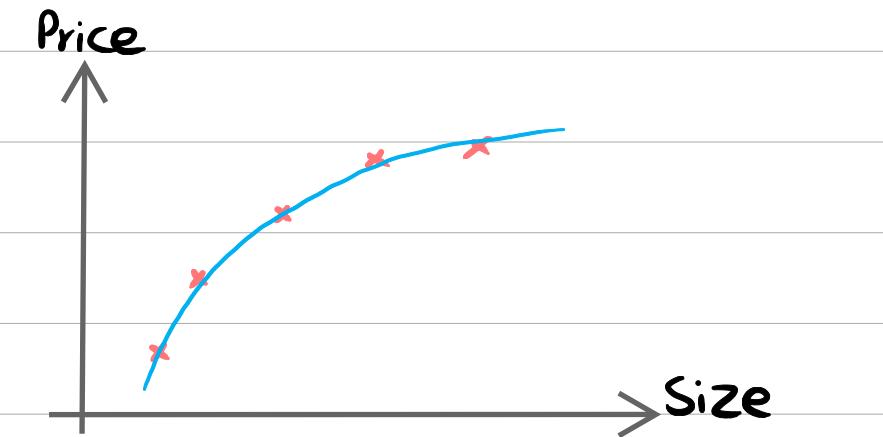
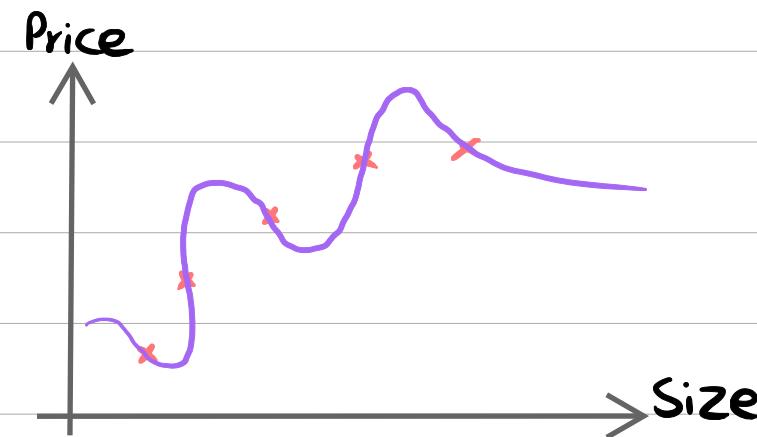
$$-0.0001x^4 + 10$$

Small values for W_j

Regularization:

بسیار کاهش خطای شود و همین تواند overfit, underfit جلوگیری کند. چه طور؟

همین تواند روش خوبی برای انتخاب دیگری باشد، چه طور؟



$$W_1x + W_2x^2 + W_3x^3 + W_4x^4 + b$$

$$W_1x + W_2x^2 + \underbrace{W_3x^3}_{\approx 0} + \underbrace{W_4x^4}_{\approx 0} + b$$

Make W_3, W_4 really small (≈ 0)

$$\min \left[\frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 \right]$$

$$\underbrace{1000 W_3^2}_1 + \underbrace{1000 W_4^2}_2$$

0.001 0.002

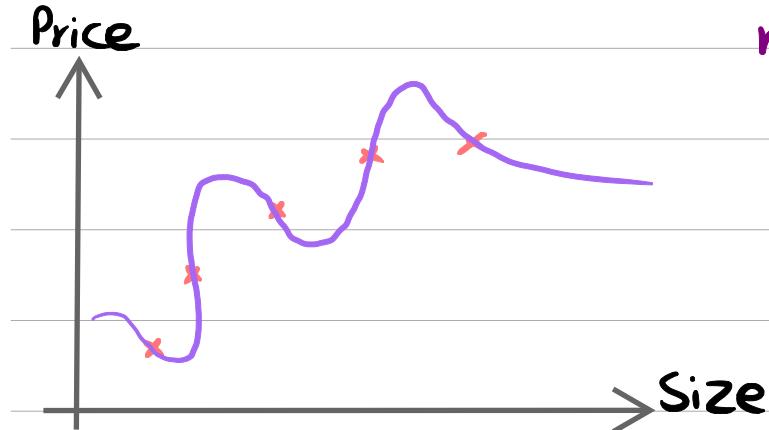
x_1	x_2	x_3	x_4	x_5	...	x_{100}	y
Size	bedroom	floors	age	avg income	...	distance coffee shop	Price

$W_1, W_2, W_3, \dots, W_{100}, b$

$n=100$

$$J(W, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n W_j^2 + \frac{\lambda}{2m} b^2$$

λ : Lambda Regularization Parameter



$$\min J(W, b) = \min \left[\frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n W_j^2 \right]$$

$W, b = ?$

Gradient Descent:

repeat

{

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w, b) \rightarrow = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

Regularizer term

}

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b) \rightarrow = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

$$\cancel{\frac{\partial}{\partial w_j}} \left[\frac{1}{2m} \sum_{i=1}^m \underbrace{(f_{w,b}(x^{(i)}) - y^{(i)})^2}_{w \cdot x^{(i)} + b} + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right] = \frac{1}{2m} \sum_{i=1}^m \left[(wx^{(i)} + b - y^{(i)}) \times 2 \times x_j^{(i)} \right] + \sum_{j=1}^n \frac{\lambda}{2m} \times 2 w_j$$

$$= \frac{1}{m} \sum_{i=1}^m \left[\underbrace{(w \cdot x^{(i)} + b - y^{(i)})}_{f_{w,b}(x)} x_j^{(i)} \right] + \sum_{i=1}^m \frac{\lambda}{m} w_j = \frac{1}{m} \sum_{i=1}^m \left[(f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j$$

$$w_j = w_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

$$w_j = w_j - \underbrace{\alpha \frac{\lambda}{m} w_j}_{w_j(1-\alpha\frac{\lambda}{m})} - \underbrace{\alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\text{usual update}}$$

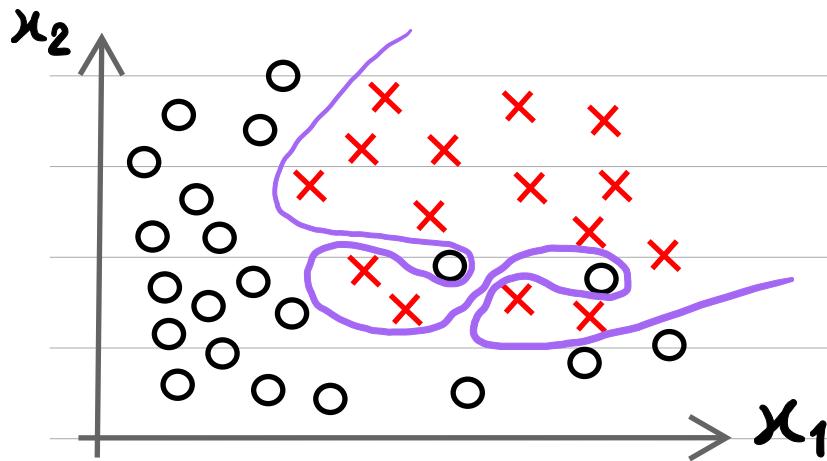
$$\hookrightarrow \alpha = 0.01, \lambda = 1 \implies \alpha \frac{\lambda}{m} = 0.01 \times \frac{1}{50} = 0.0002$$

m = 50

$$w_j \overset{0.9998}{\cancel{(1-0.0002)}}$$

Classification \rightarrow Logistic Regression

Overfit \leadsto high variance \uparrow



$$\begin{aligned} Z = & W_1 x_1 + W_2 x_2 + W_3 x_1^2 x_2 + W_4 x_1^2 x_2^2 \\ & + W_5 x_1^2 x_2^3 + W_6 x_1^3 x_2 + \dots + b \end{aligned}$$

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{w,b}(x^{(i)})) + (1-y^{(i)}) \log(1-f_{w,b}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

$\min J(w, b) \leadsto w_j \downarrow$

GD:

repeat

{

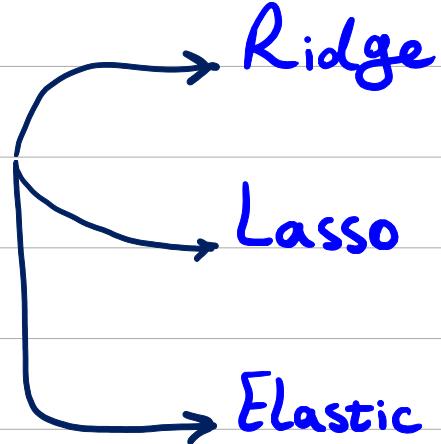
$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w, b) \rightarrow = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

Logistic Regression $\rightarrow f_{w,b}(x) = \frac{1}{1+e^{-wx+b}}$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b) \rightarrow = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

}

Regularization



$$L2 = \sum_{j=0}^n w_j^2$$

alpha $\rightarrow \alpha$

$$L1 = \sum_{j=0}^n |w_j|$$

L1, L2

Ridge \rightarrow مکان ضرایب کوچک ہی شوند

$\alpha = 0 \rightarrow$ Ridge X

$$\text{Loss} = \sum_{i=1}^m (f_{w,b}(x) - y^{(i)})^2 + \lambda \sum_{j=0}^n w_j^2$$

$\alpha = \infty \rightarrow$ مکان ضرایب کوچک و نزدیک ب صفر ہی شوند.

$\alpha \uparrow$

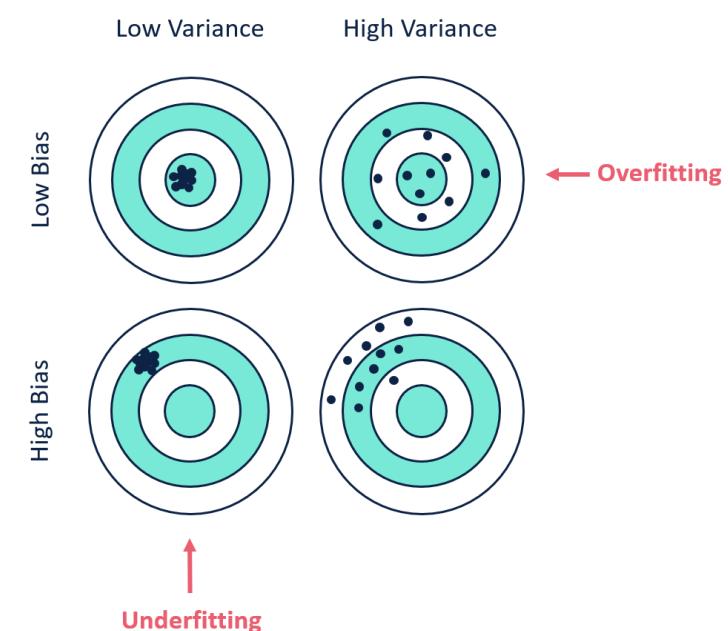
bias \uparrow

Variance \downarrow

$\alpha \downarrow$

bias \downarrow

Variance \uparrow



هنرایب را که جگ های لند و برقی از هنرایب را به ۵ هی رساند.

$$\text{Loss} = \sum_{i=1}^m \left(f_{w,b}(x_i) - y^{(i)} \right)^2 + \lambda \sum_{j=0}^n |w_j|$$

$\alpha = 0 \rightarrow$ همچنانچه وزنی حذف نمی شود.

$\alpha = \infty \rightarrow$ تمام وزنها را حذف می کند.

Elastic Net:

$\hookrightarrow L1, L2$

ترکیبی از Lasso, Ridge

$$\text{Loss} = \frac{1}{2m} \sum_{i=1}^m \left(f_{w,b}(x_i) - y_i \right)^2 + (1-\lambda) \times \frac{\alpha}{2} \times \underbrace{\sum_{j=1}^n w_j^2}_{L2} + \lambda \alpha \underbrace{|w_j|}_{L1}$$

$L1_rate = 1 \rightarrow L1 \rightarrow \text{lasso}$

$L1_rate = 0 \rightarrow L2 \rightarrow \text{Ridge}$

x	y
1	2
2	3
3	5

$$\hat{y} = wx + b$$

$$\text{Loss} = \text{MSE} = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

$$w=1.5, b=0.5$$

$$\text{Loss} = \text{Reg. MSE} = \text{MSE} + \lambda \sum_{j=1}^m w_j^2$$

$$\lambda = 0.1$$

$$\hat{y}_1 = wx + b = 1.5 \times 1 + 0.5 = 2$$

$$\hat{y}_2 = 1.5 \times 2 + 0.5 = 3.5$$

$$\hat{y}_3 = 1.5 \times 3 + 0.5 = 5$$

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2 = \frac{1}{3} [(2-2)^2 + (3-3.5)^2 + (5-5)^2] = \frac{1}{3} [0 + 0.25 + 0] = 0.083$$

$$\text{Reg. MSE} = \text{MSE} + \lambda w^2 = 0.083 + (0.1 \times (1.5)^2) = 0.308$$

$$\text{Reg. MSE} > \text{MSE}$$

برای کاهش MSE چه باید کرد؟

اگر w را کاهش دهیم MSE کاهشی باید و از طرفی با کاهش w مدل ازبیش برآورده نجات یابد.

یک تعادل ایکاری لندین fit-well و دهنای کوچک که تواند باست Generalization بختر داده های دیده نشده شود. * Reg.

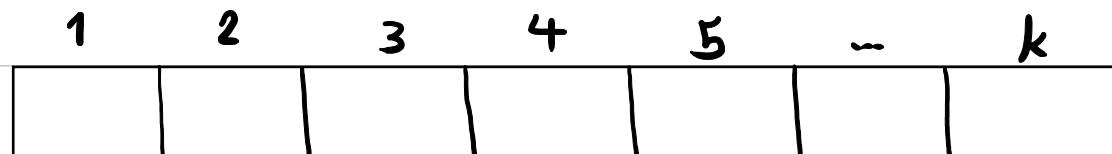
Cross Validation

۸. آیا روشی برای پیدا کردن هایپر پارامتر (α)

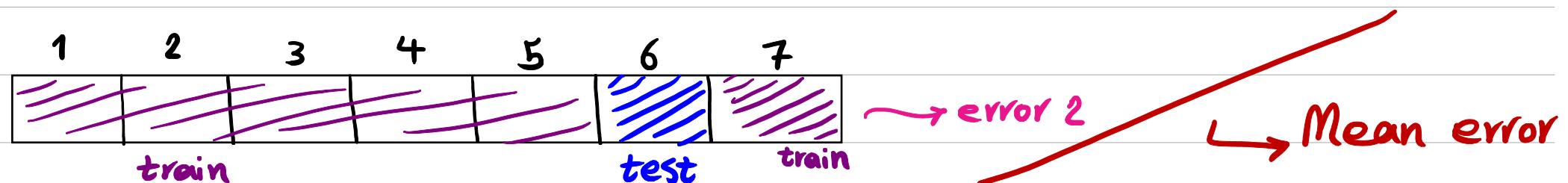
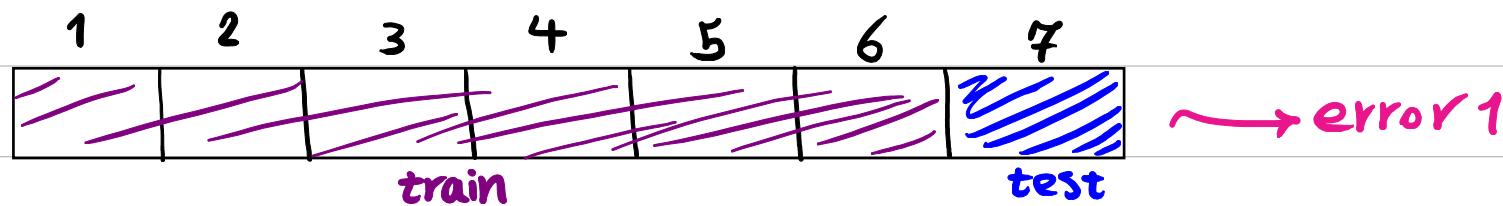
	x_1	x_2	x_3	x_4	x_5	y
train						
test						

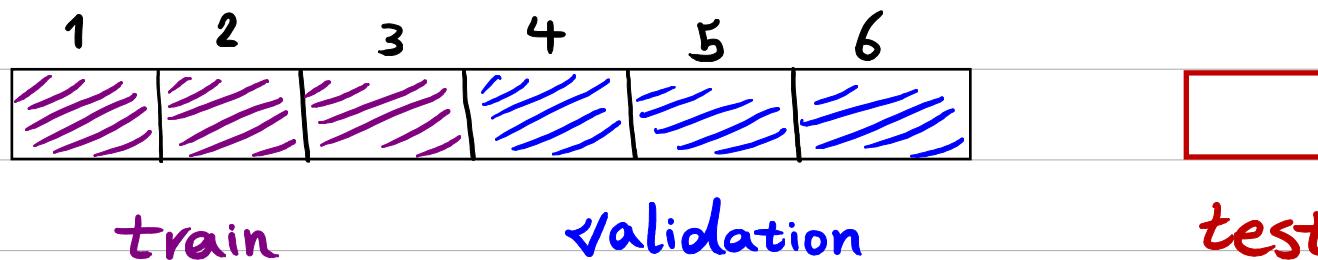
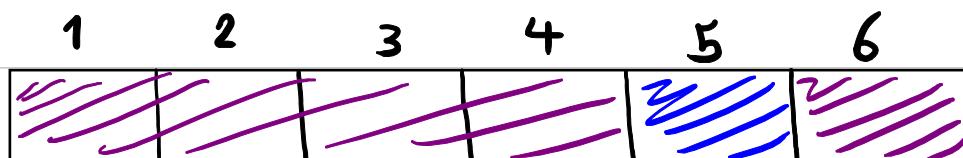
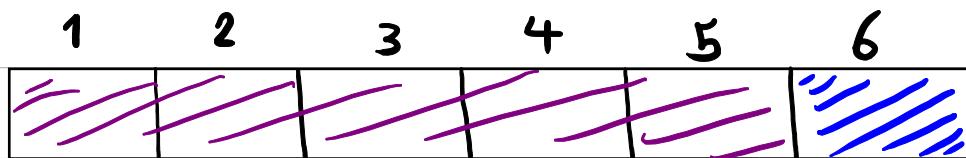
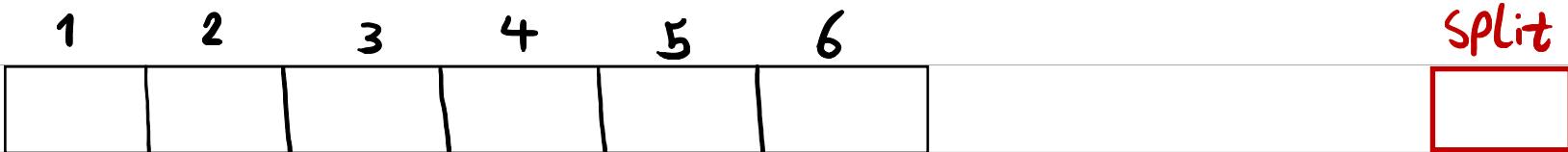
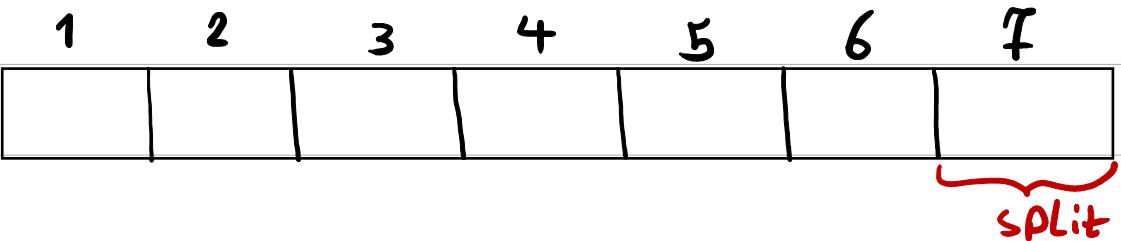
نمایی ملت؟

k -fold:



$k=7$, $CV=2$





x

y

1	2
2	4
3	6
4	8
5	10

k-fold cross-validation

①

1	2	3	4	5
---	---	---	---	---

$$\rightarrow \text{MSE} = 0$$



$$\begin{aligned} x &= [1] \\ y &= [2] \end{aligned}$$

$$\text{train} \rightarrow x = [2, 3, 4, 5]$$

$$y = [4, 6, 8, 10]$$

②

1	2	3	4	5
///	///	///	///	///

$$\rightarrow \text{MSE} \simeq 7.9 \times 10^{-31}$$

③

1	2	3	4	5
///	///	///	///	///

$$\rightarrow \text{MSE} = 0$$

④

1	2	3	4	5
///	///	///	///	///

$$\rightarrow \text{MSE} = 0$$

⑤

1	2	3	4	5
///	///	///	///	///

$$\rightarrow \text{MSE} = 0$$

$$\text{avg_mse} = \frac{1}{m} \sum^k \text{MSE} = 1.5 \times 10^{-31}$$



