

# Linear Regression

Lecture 3

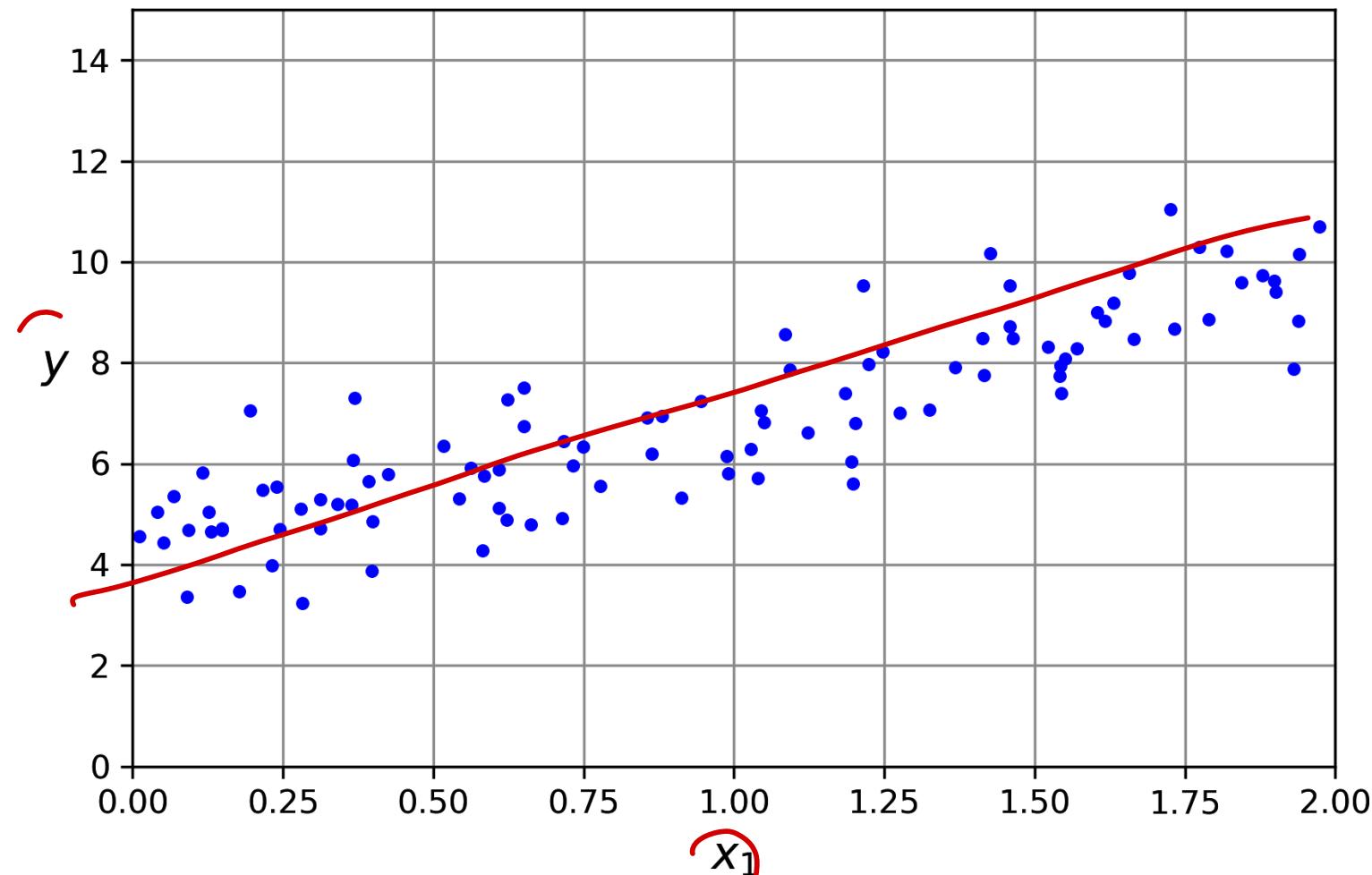
# Linear Regression

- Generating models which are a linear function of the input features.

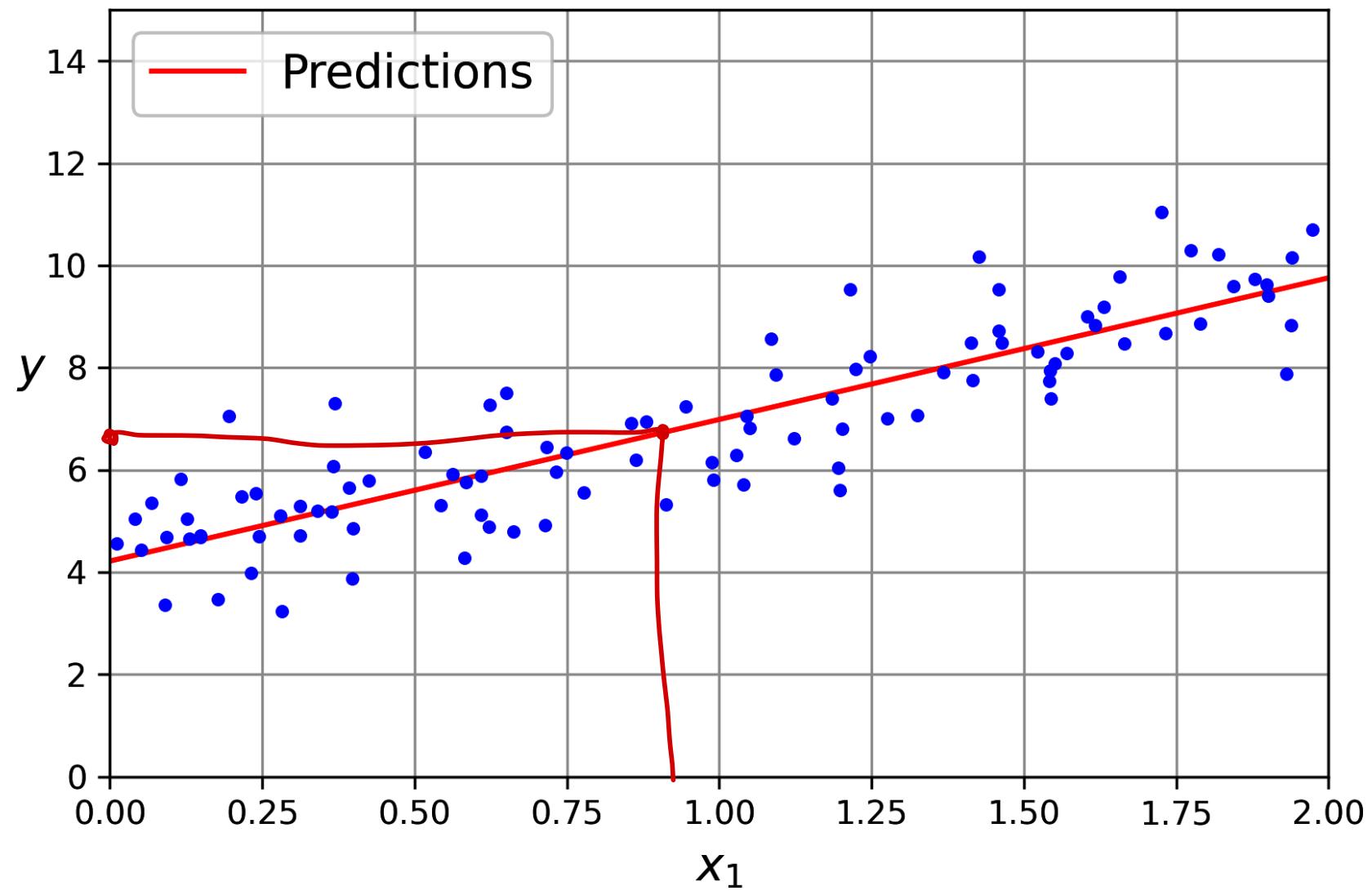
ML

D  
Pattern  
Machine

$y = m_1x + c$



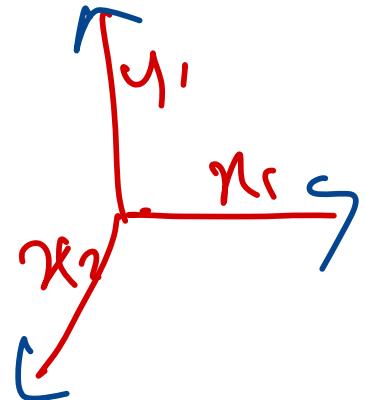
# Linear Regression



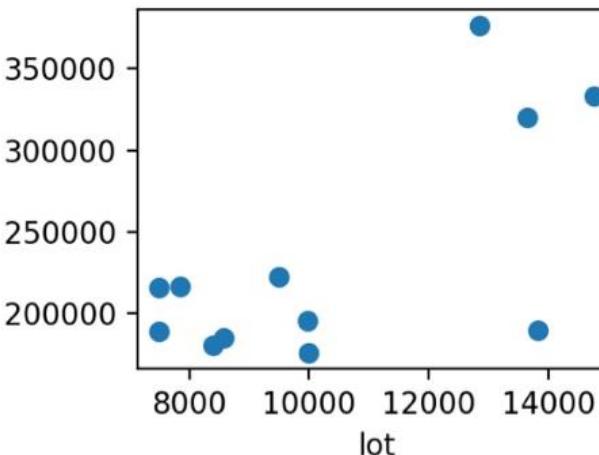
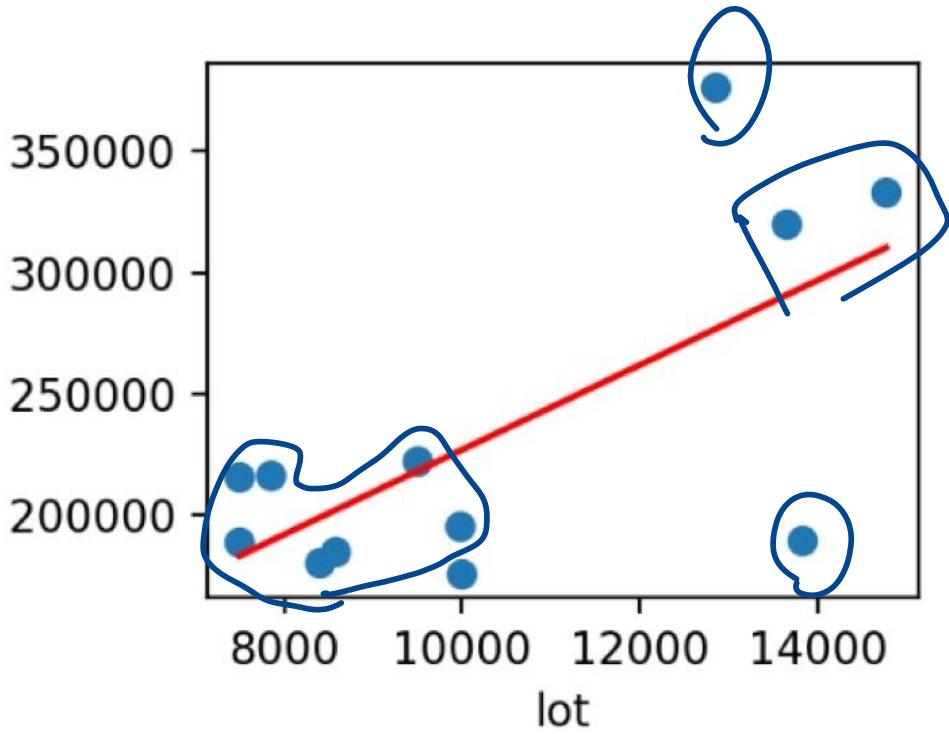
# Hypothesis

CS229

- ▶ A **hypothesis** or a prediction function is function  $h : \mathcal{X} \rightarrow \mathcal{Y}$ 
  - ▶  $\mathcal{X}$  is an image, and  $\mathcal{Y}$  contains "cat" or "not."
  - ▶  $\mathcal{X}$  is a text snippet, and  $\mathcal{Y}$  contains "hate speech" or "not."
  - ▶  $\mathcal{X}$  is house data, and  $\mathcal{Y}$  could be the price.
- ▶ A **training set** is a set of pairs  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$   
s.t.  $x^{(i)} \in \mathcal{X}$  and  $y^{(i)} \in \mathcal{Y}$  for  $i = 1, \dots, n$ .
- ▶ Given a training set our goal is to produce a *good* prediction function  $h$ 
  - ▶ Defining "*good*" will take us a bit. It's a modeling question!
  - ▶ We will want to use  $h$  on new data not in the training set.
- ▶ If  $\mathcal{Y}$  is continuous, then called a regression problem.
- ▶ If  $\mathcal{Y}$  is discrete, then called a classification problem.



# Housing Data



	SalePrice	Lot.Area
4	189900	13830
5	195500	9978
9	189000	7500
10	175900	10000
12	180400	8402
22	216000	7500
36	376162	12858
47	320000	13650
55	216500	7851
56	185088	8577
58	222500	9505
59	333168	14774

# Hypothesis + More Features

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

Parameters.

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

$x_1$   
Living area (feet<sup>2</sup>)  
2104.  
1600.  
2400.

$x_2$   
#bedrooms  
3  
3  
3

$y^i$   
Price (1000\$)  
400  
330  
369

Two features – size, no of bedrooms ...

Target variable is price.

$$h(x) = \sum_{j=0}^n \theta_j x_j$$

$h \Rightarrow$  hypothesis

$\theta \Rightarrow$  parameters

$n \Rightarrow$  no. of features.

$x \Rightarrow$  feature / input

$M \Rightarrow$  no. of training ex

$y \Rightarrow$  output / target var.

$(x_i, y_i)$  training example

$(x_i^i, y_i^i)$  example no.  
 $i^{th}$  example

$$400 = \theta_0 + 2104 \theta_1 + 3 \theta_2$$

$$330 = \theta_0 + 1600 \theta_1 + 3 \theta_2$$

$$369 = \theta_0 + 2400 \theta_1 + 3 \theta_2$$

$\left( \begin{array}{c} 1 \\ 2104 \\ 3 \end{array} \right)$

$\left( \begin{array}{c} 1 \\ 1600 \\ 3 \end{array} \right)$

$\left( \begin{array}{c} 1 \\ 2400 \\ 3 \end{array} \right)$

## Cost Function

$$h_{\theta}(x) \Rightarrow y$$

Find  $\underline{\theta}$  to  $\min J(\theta)$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

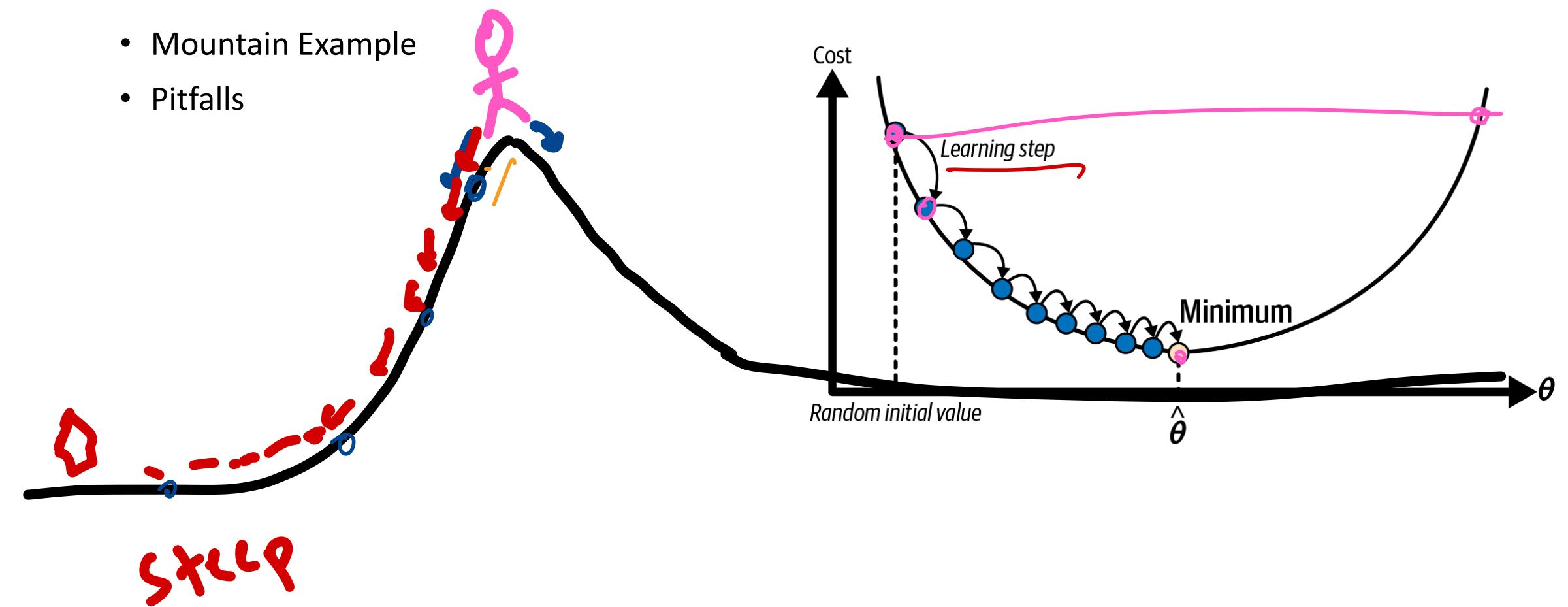
Ordinary least square

MSE

Convex function

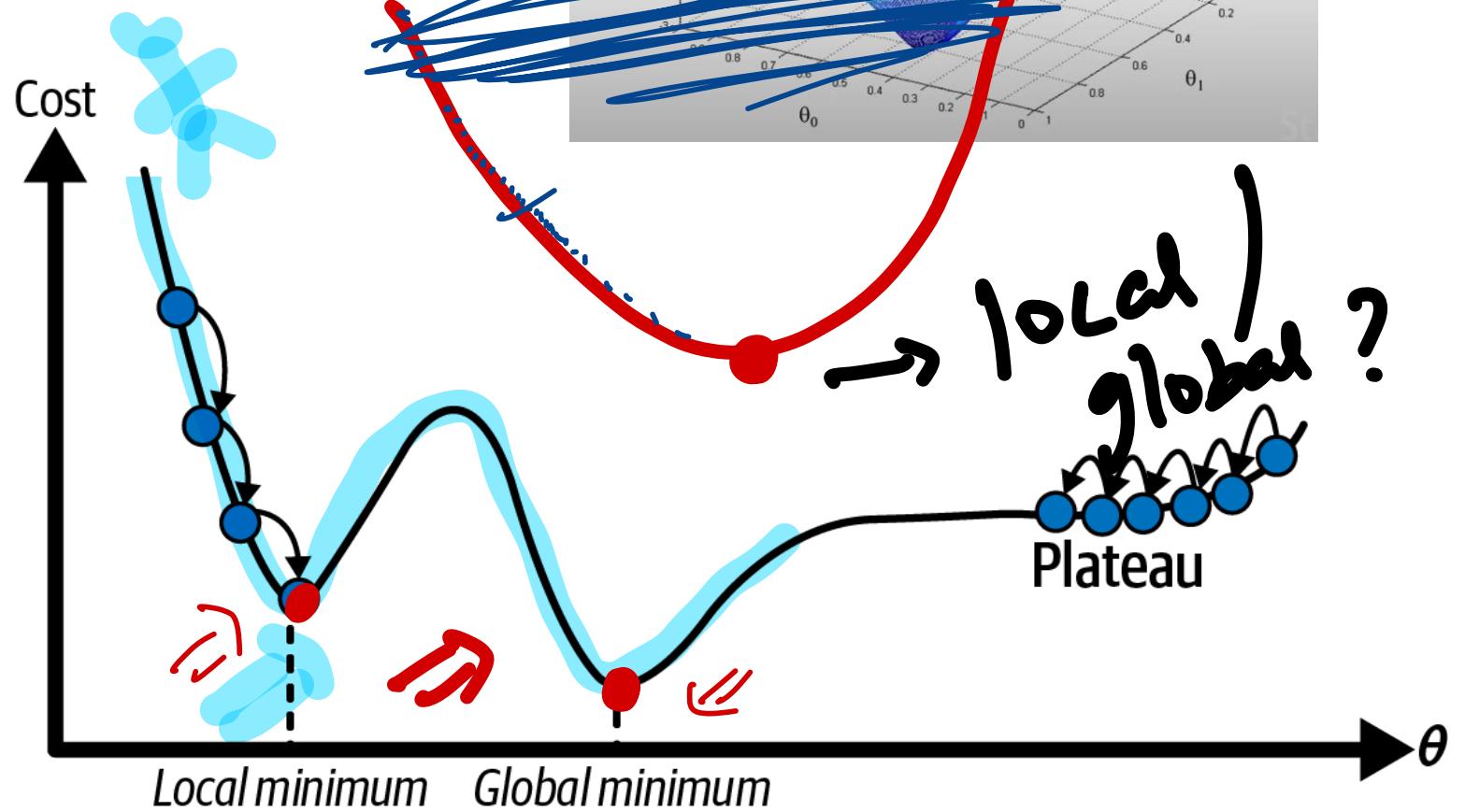
# Least Squares Optimization Problem

- *Gradient descent* is a generic optimization algorithm capable of finding optimal solutions to a wide range of problems. The general idea of gradient descent is to tweak parameters iteratively in order to minimize a cost function.
- Mountain Example
- Pitfalls



# Pitfalls

- Local optima for irregular terrain
- The MSE cost function is convex which means that there is no local minima
- Convex function any two points on the curve, the line segment joining them is never below the curve.



# Iterations for finding $\theta$

start with some value of  $\theta$   
 $\theta = 0$   
 Find  $\theta_0, \theta_1, \theta_2, \dots$  s.t. minimise  $J(\theta)$

$\theta_j^{t+1} = \theta_j^t - \frac{\alpha}{\delta \theta} (\frac{\partial J(\theta)}{\partial \theta_j})$   
 , learning rate

$$\begin{aligned}\frac{\partial}{\partial \theta_j} (J(\theta)) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} \sum (h_\theta(x) - y)^2 \\ &= \cancel{2x} \frac{1}{2} (h_\theta(x) - y) \frac{\partial}{\partial \theta_j} (h_\theta(x) - y)\end{aligned}$$

$\Rightarrow$

$\alpha \Rightarrow$  learning rate.

$0 - 10^0$

$0 - 5^0$

$2^0 - 3^0$

$0 - 10000 \text{ UD}$

-1 - 1

$\begin{cases} v.\text{large} \\ 1 \end{cases} n_1 + \underline{v.\text{small}}$

$$(h_\theta(x) - y) \left[ \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n - y) \right]$$

All terms are 0 except for terms corresponding to  $j$

$$\theta_j^{t+1} = \theta_j^t - \alpha \left( \frac{\partial J(\theta)}{\partial \theta_j} \right)$$

$$\theta_j^{t+1} = \theta_j^t + \alpha (y - h_\theta(x)) x_j^{(i)}$$

Repeat until convergence for  $j = (0, 1, \dots, n)$

# Iterations for finding $\Theta$

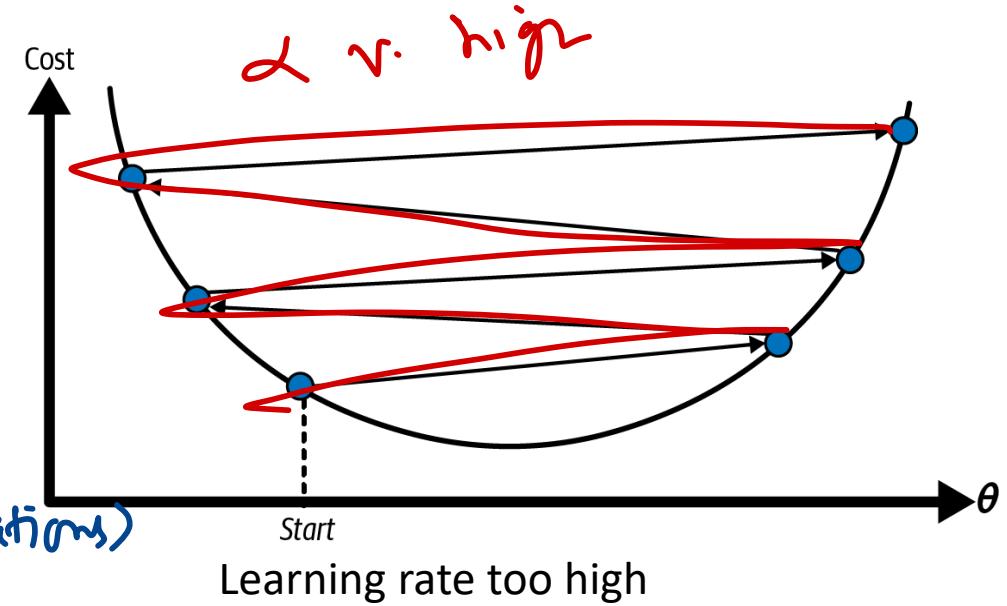
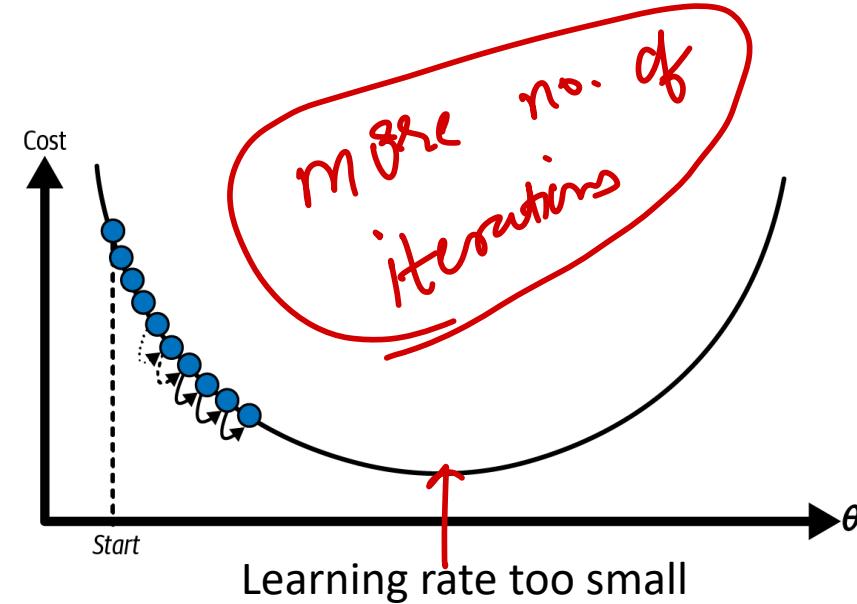
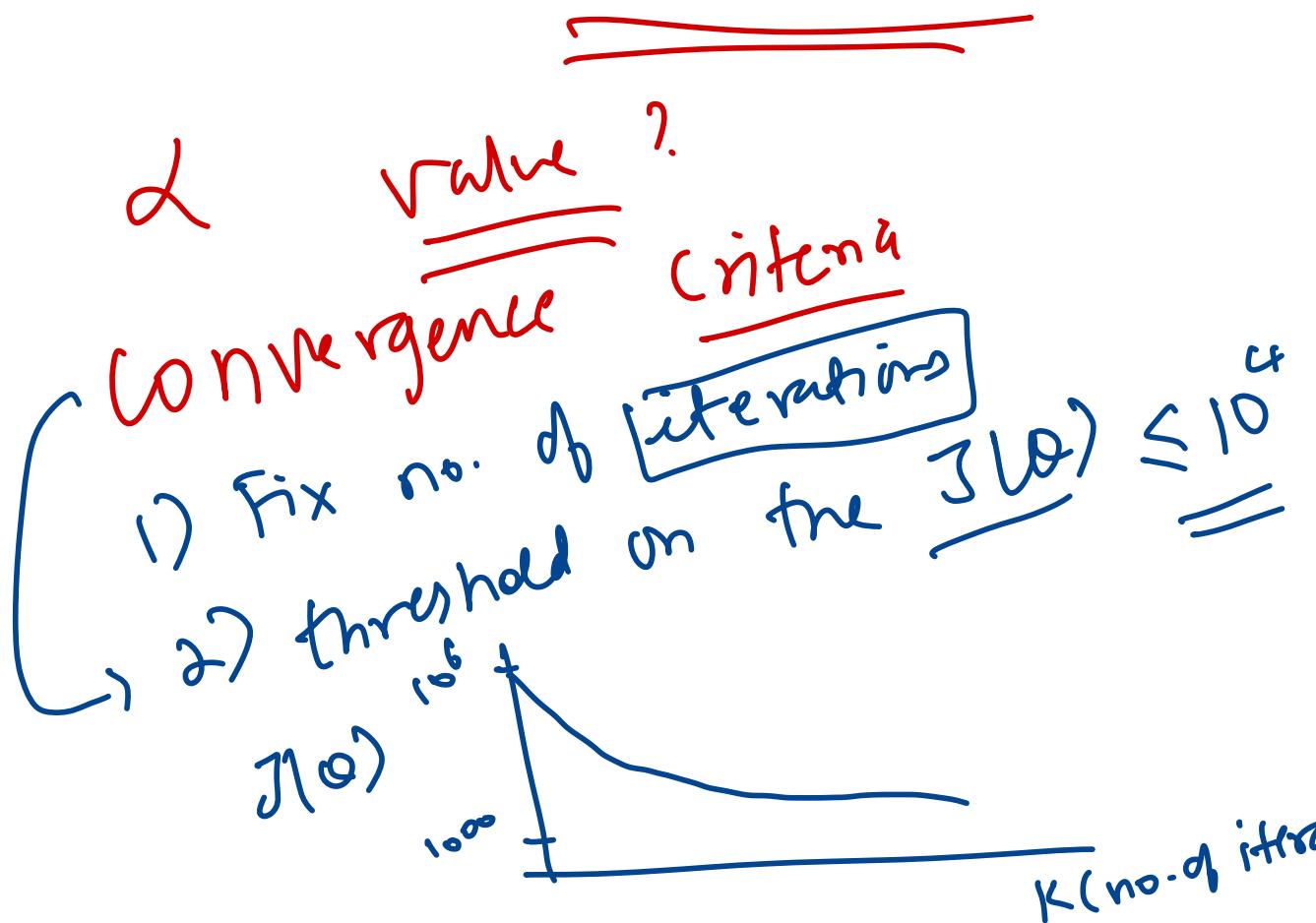
Prescritic no. of iterations

$\Rightarrow$  1000

$$J < 10^{-6}$$

# Learning Rate

- Learning rate hyperparameter

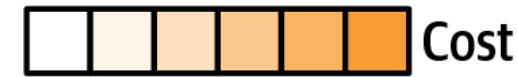


# Feature Scaling for Gradient Descent

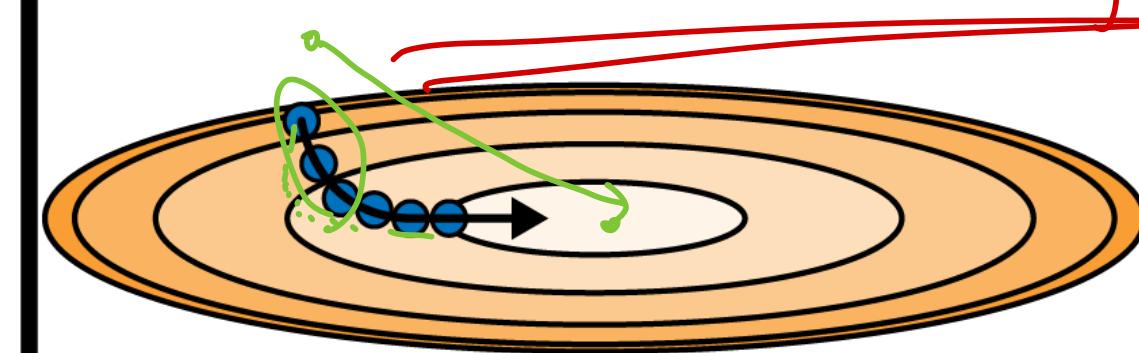
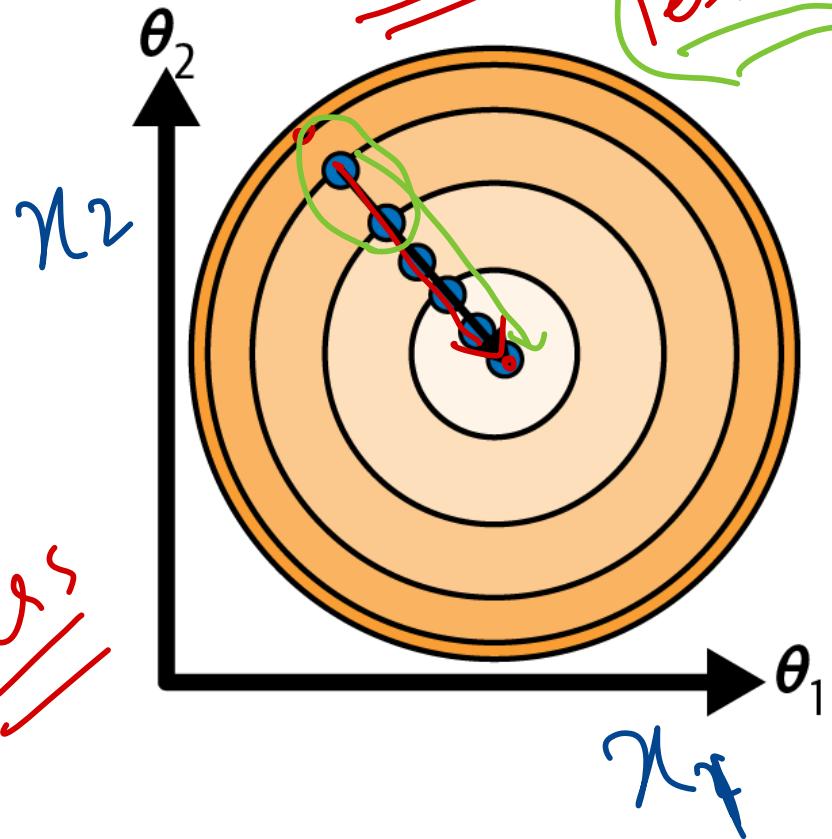
49x3

quicker  
less time  
@ higher cost

$\theta_0$   
 $\theta_1$   
 $\theta_2$

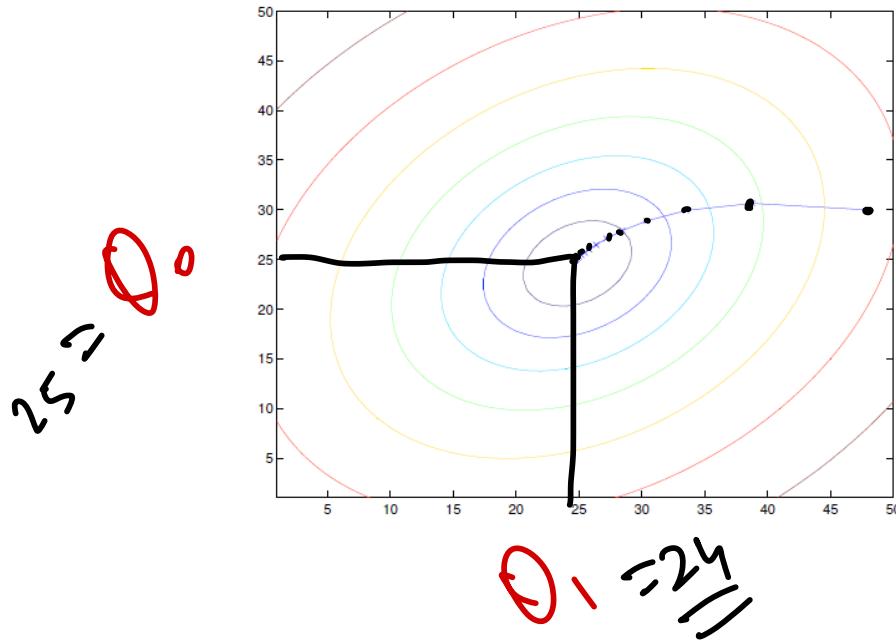


more fine @ higher cost



(more convergence time)

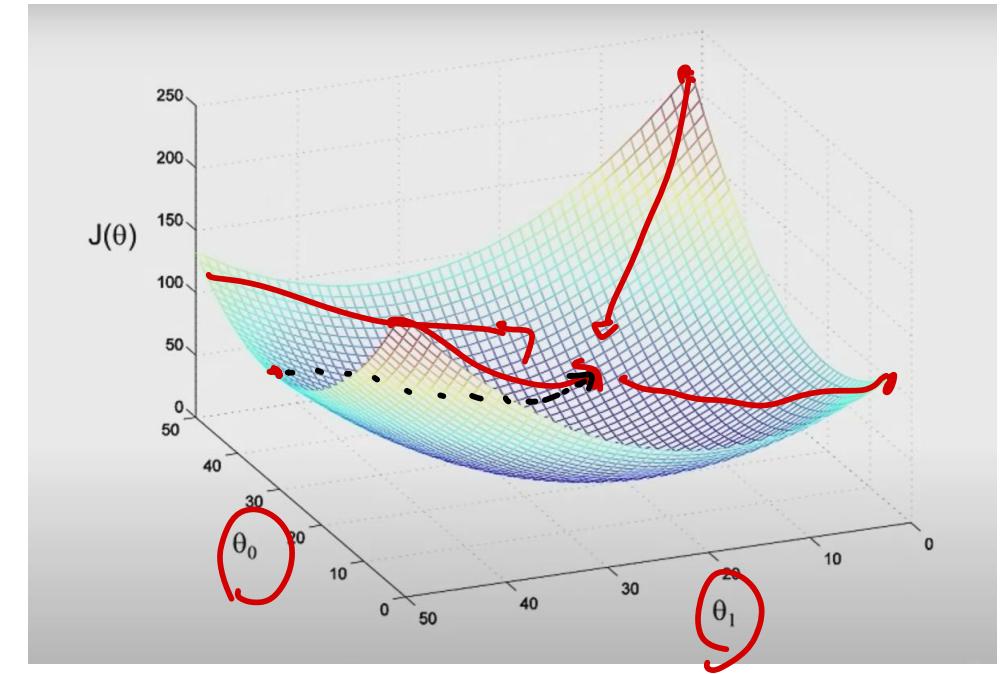
# Cost Function



Convex Function  
Slices of the convex function  
Contours of the quadratic function.

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

$J(\theta)$



# Solution

$$n=2 \quad m=49$$

$$\theta_0 = 89.60, \theta_1 = 0.1392, \theta_2 = -8.738$$

- ~~100~~ – housing price data
- Parameters are updated according to the gradient of the error with respect to that single training example only

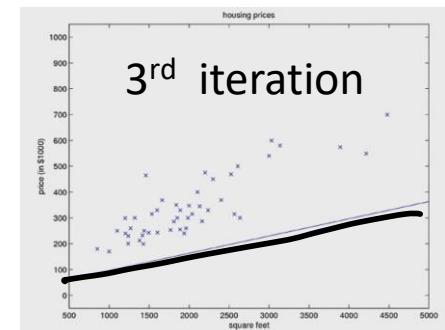
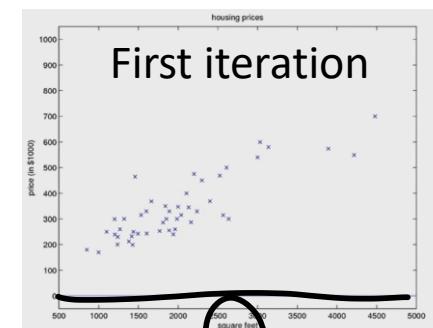
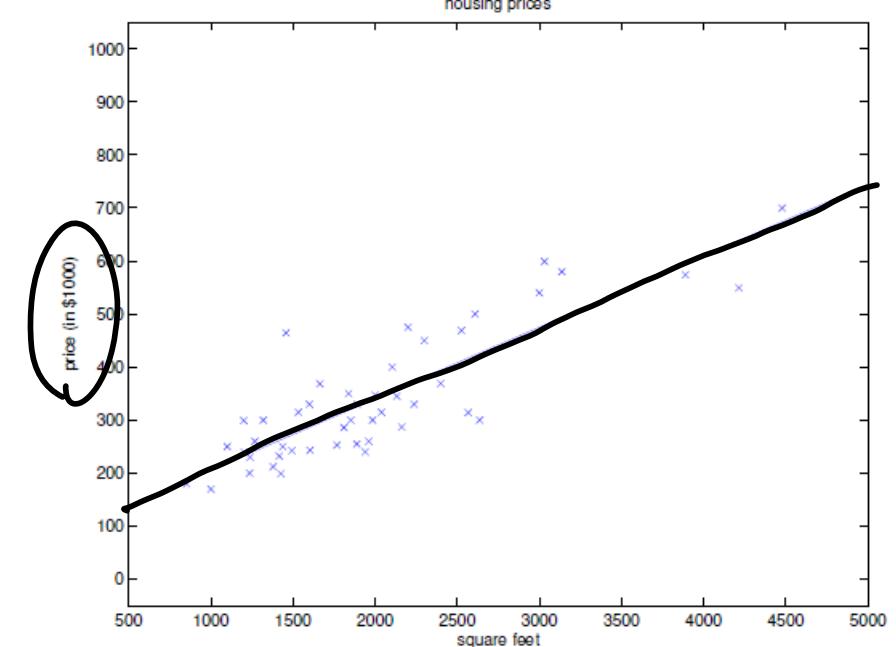
Loop {

    for  $i = 1$  to  $n$ , {

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}, \quad (\text{for every } j)$$

}

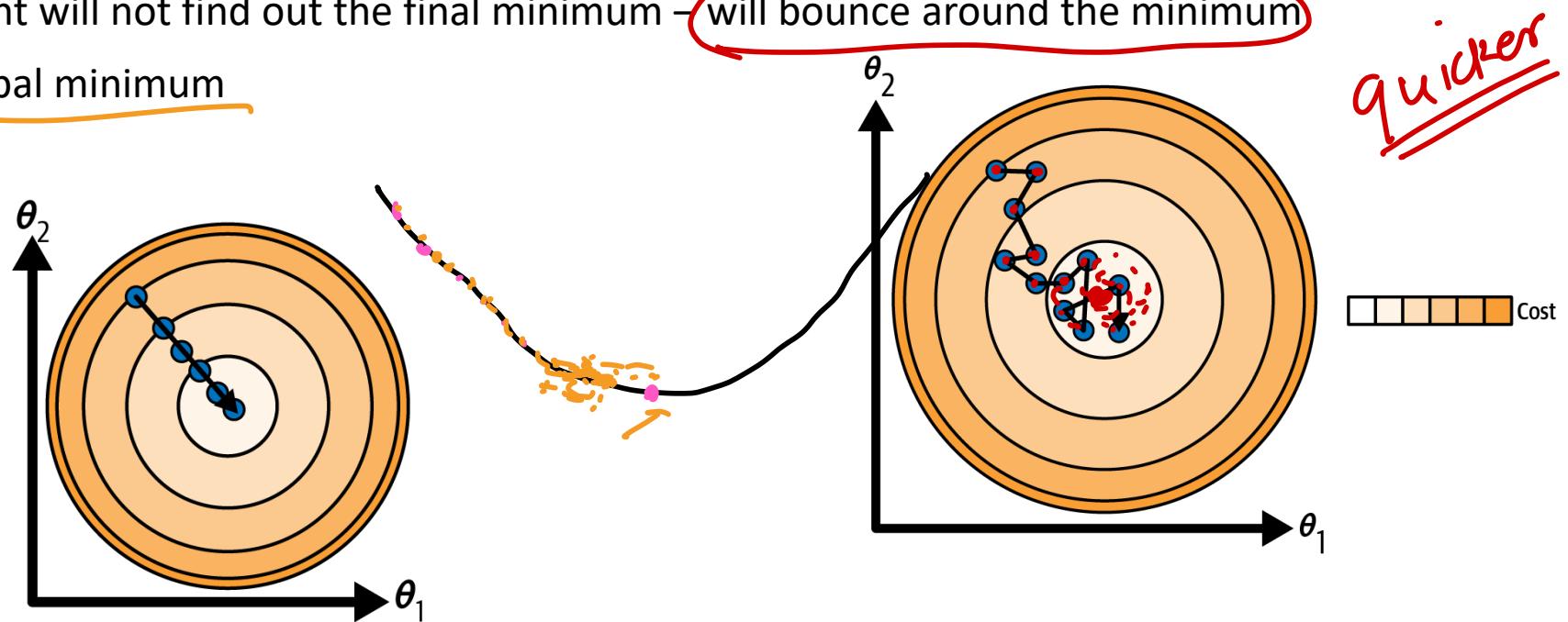
}



# Batch Gradient Descent vs. Stochastic GD.

- The main problem with batch gradient descent is the fact that it uses the whole training set to compute the gradients at every step, which makes it very slow when the training set is large.  
*slows down*
- Stochastic gradient descent* picks a random instance in the training set at every step and computes the gradients based only on that single instance.
- Stochastic gradient descent will not find out the final minimum – will bounce around the minimum
- Will oscillate near the global minimum
- ~~Adaptive learning rate~~

*reduce*



# Stochastic Gradient Descent with Dynamic Learning Schedule

```
n_epochs = 50
t0, t1 = 5, 50 # learning schedule hyperparameters

def learning_schedule(t):
    return t0 / (t + t1)

np.random.seed(42)
theta = np.random.randn(2, 1) # random initialization

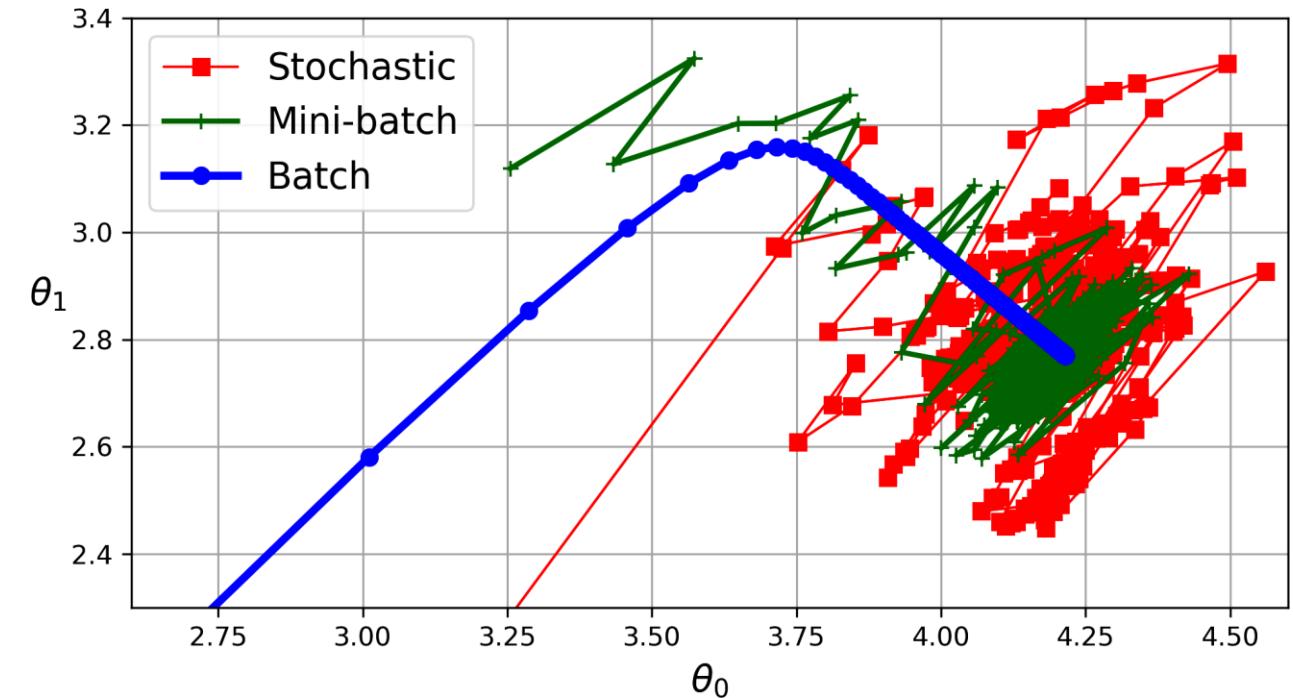
for epoch in range(n_epochs):
    for iteration in range(m):
        random_index = np.random.randint(m)
        xi = X_b[random_index : random_index + 1]
        yi = y[random_index : random_index + 1]
        gradients = 2 * xi.T @ (xi @ theta - yi) # for SGD, do not divide by m
        eta = learning_schedule(epoch * m + iteration)
        theta = theta - eta * gradients
```

dynamic  $\eta$

Linear

# Mini Batch GD

- Computes gradients on small random sets of instances called mini-batches.
- A performance boost is obtained in terms of hardware optimization of matrix operations, esp. for GPUs.
- Algorithm is less erratic as compared to stochastic GD.
- Batch GD finds the exact solution – takes more time though.



# Normal Equation Derivation

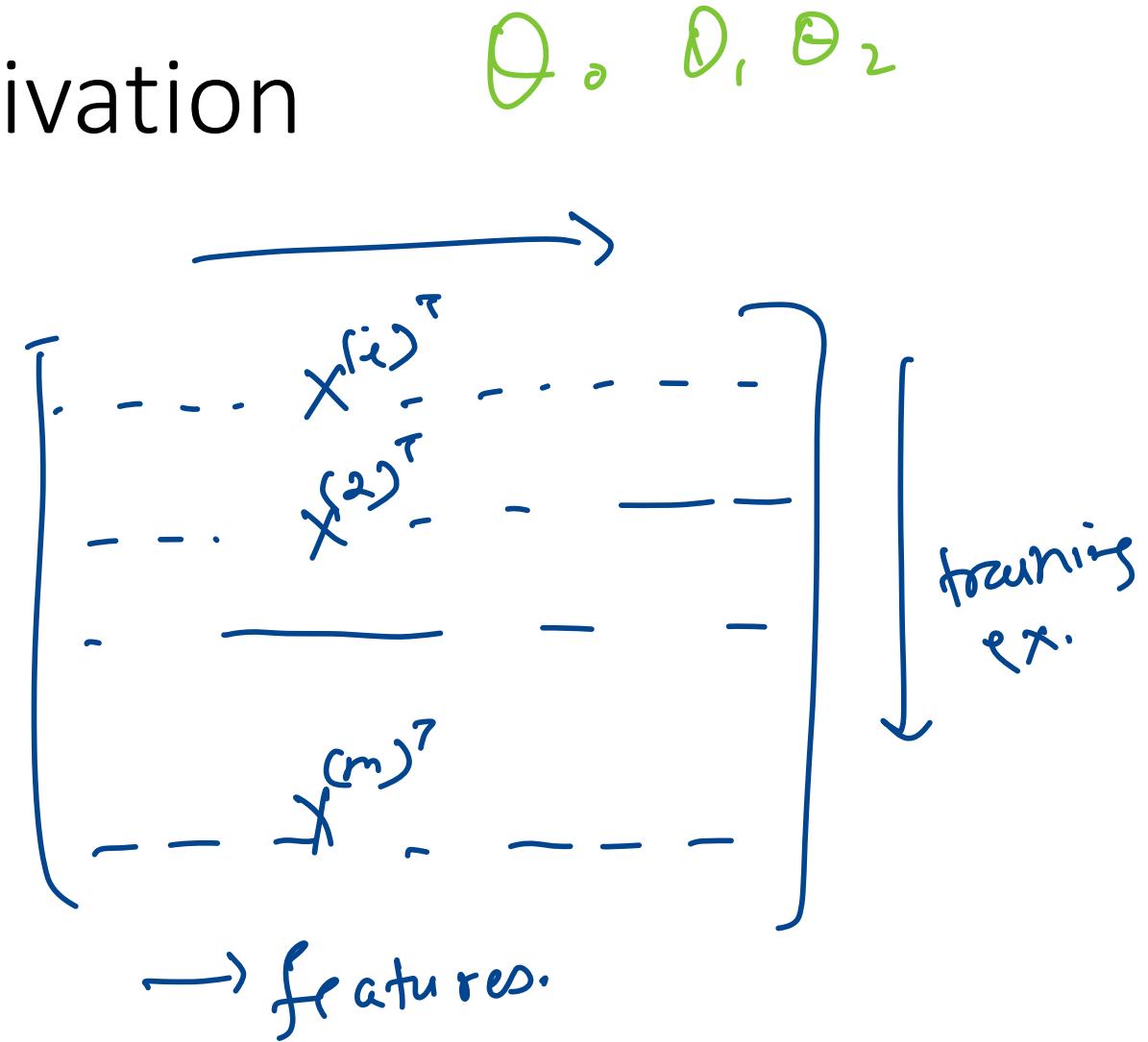
$$\hat{\theta} = (X^T X)^{-1} X^T y \Leftarrow$$

$$\theta_0 = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

$n+1 \times 1$

$$x^i = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

$m \times (n+1)$



$$X \theta \Rightarrow \boxed{m \times 1}$$

$$x_\theta = \begin{bmatrix} x^\top \theta \\ \vdots \\ x^m \theta \end{bmatrix} = \begin{bmatrix} h_\theta(x^1) \\ \vdots \\ h_\theta(x^m) \end{bmatrix}$$

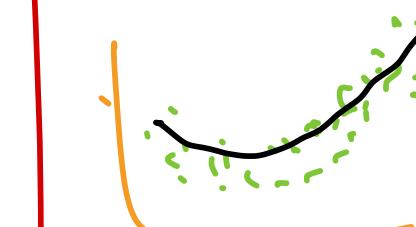
$$\nabla_\theta J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix}, \quad \theta \in \mathbb{R}^{n+1}$$

(n+1 x 1)

!  $\frac{\partial J}{\partial \theta_n}$

$$\hat{y} = \begin{bmatrix} y^1 \\ \vdots \\ y^m \end{bmatrix}$$

$$(x_\theta - y) = \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}$$



$$\frac{1}{2} (x\theta - y)^T (x\theta - y) \rightarrow J(\theta)$$

$$J^2 = \sum_i z_i^2$$

$$\frac{1}{2} \sum (h_\theta(x^{(i)} - y^{(i)})^2$$

# Normal Equation

- Gives the result directly for the MSE equation.

gradient

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (\underline{X\theta} - \vec{y})^T (\underline{X\theta} - \vec{y}) \\ &= \frac{1}{2} \nabla_{\theta} ((X\theta)^T X\theta - (X\theta)^T \vec{y} - \vec{y}^T (X\theta) + \vec{y}^T \vec{y}) \\ &= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X)\theta - \vec{y}^T (X\theta) - \vec{y}^T (X\theta)) \\ &= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X)\theta - 2(X^T \vec{y})^T \theta) \\ &= \frac{1}{2} (2X^T X\theta - 2X^T \vec{y}) \\ &= \boxed{X^T X\theta - X^T \vec{y}}\end{aligned}$$

Thus, the value of  $\theta$  that minimizes  $J(\theta)$  is given in closed form by the equation

$$\theta = \boxed{\underline{(X^T X)}^{-1} X^T \vec{y}}$$

Large dataset

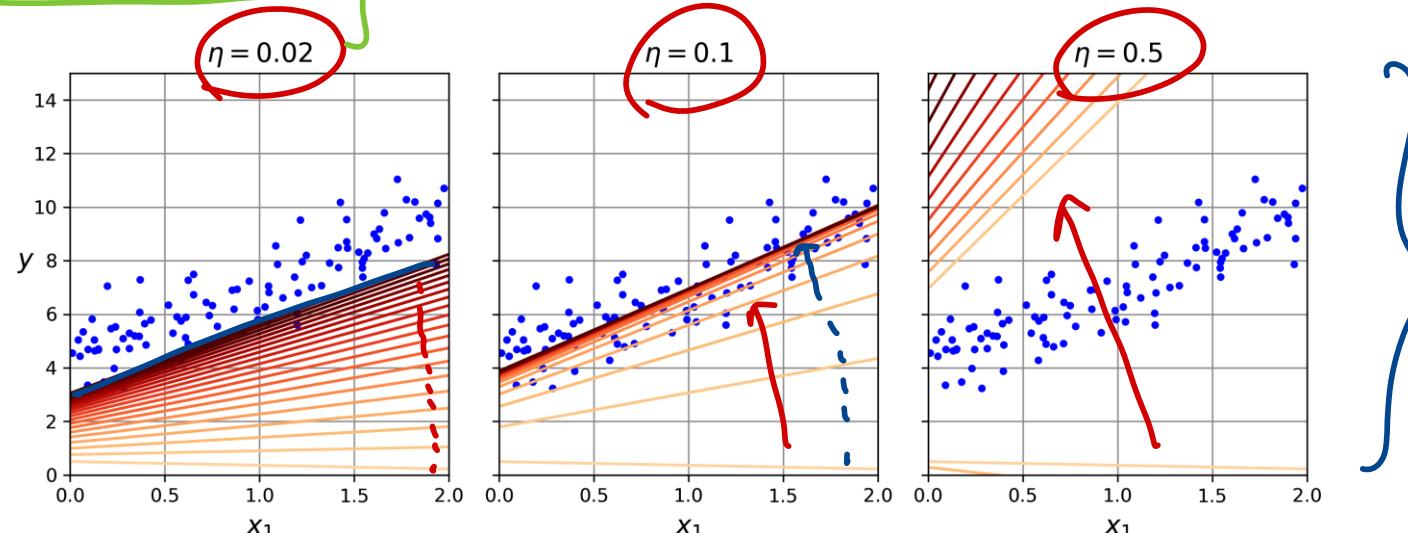
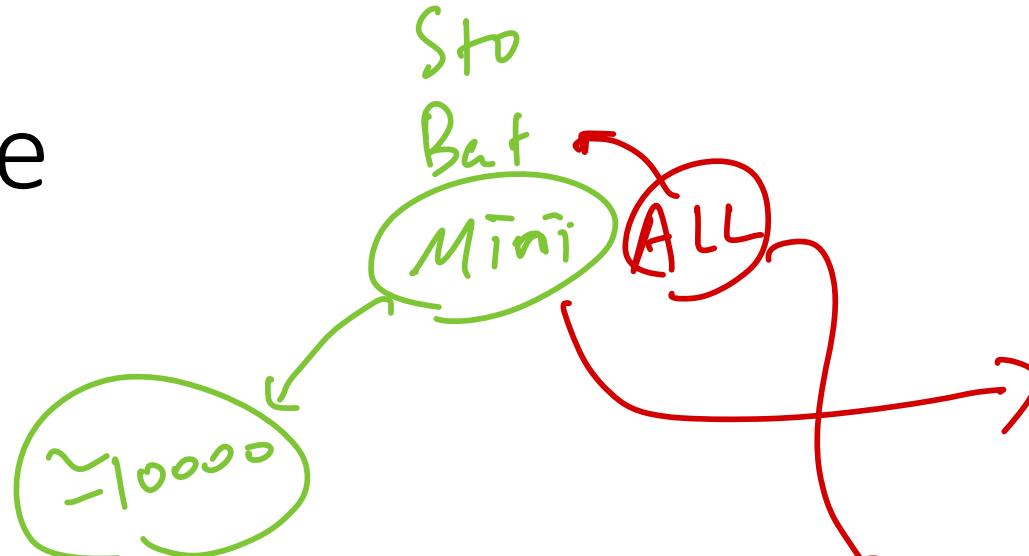
# Normal Equation code

```
eta = 0.1 # learning rate
n_epochs = 1000
m = len(X_b) # number of instances

np.random.seed(42)
theta = np.random.randn(2, 1) # randomly initialized model parameters

for epoch in range(n_epochs):
    gradients = 2 / m * X_b.T @ (X_b @ theta - y)
    theta = theta - eta * gradients

>>> theta
array([[4.21509616],
       [2.77011339]])
```

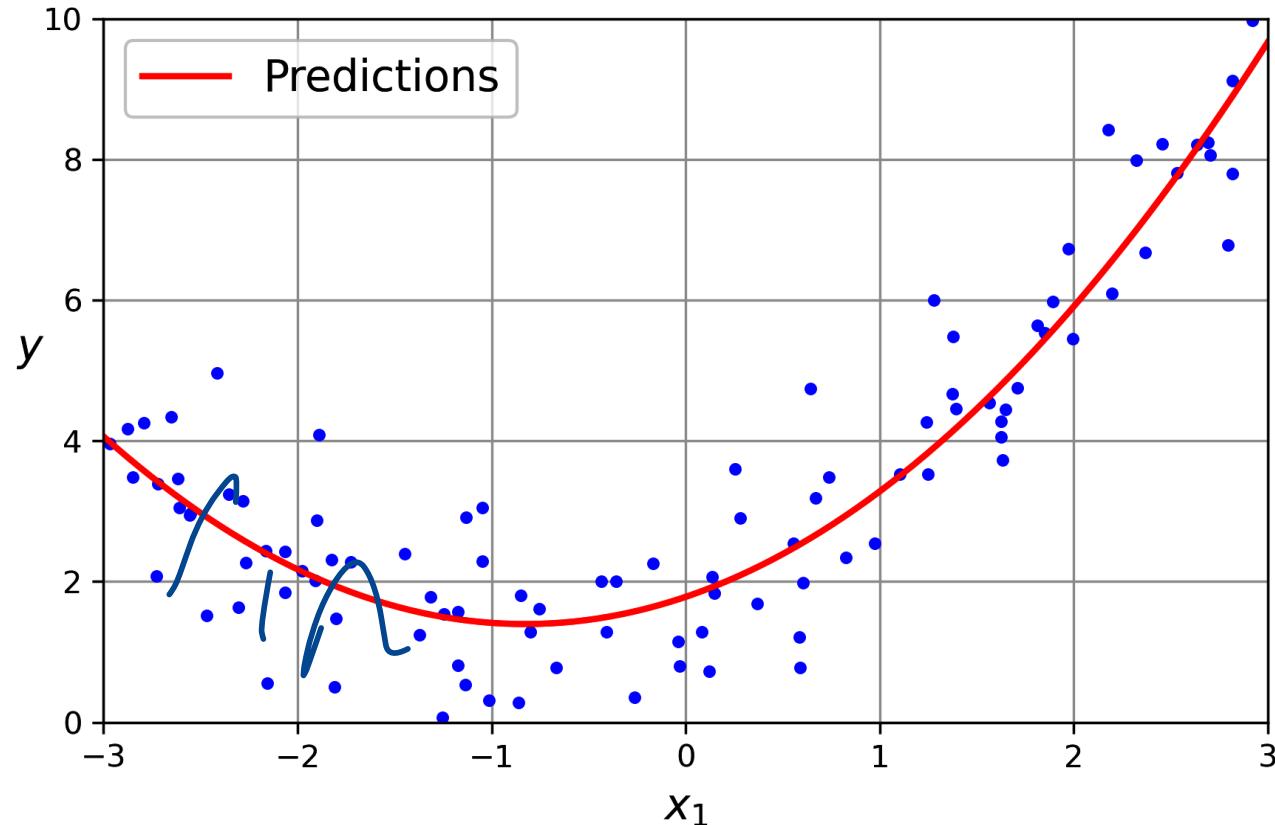


# Polynomial Regression

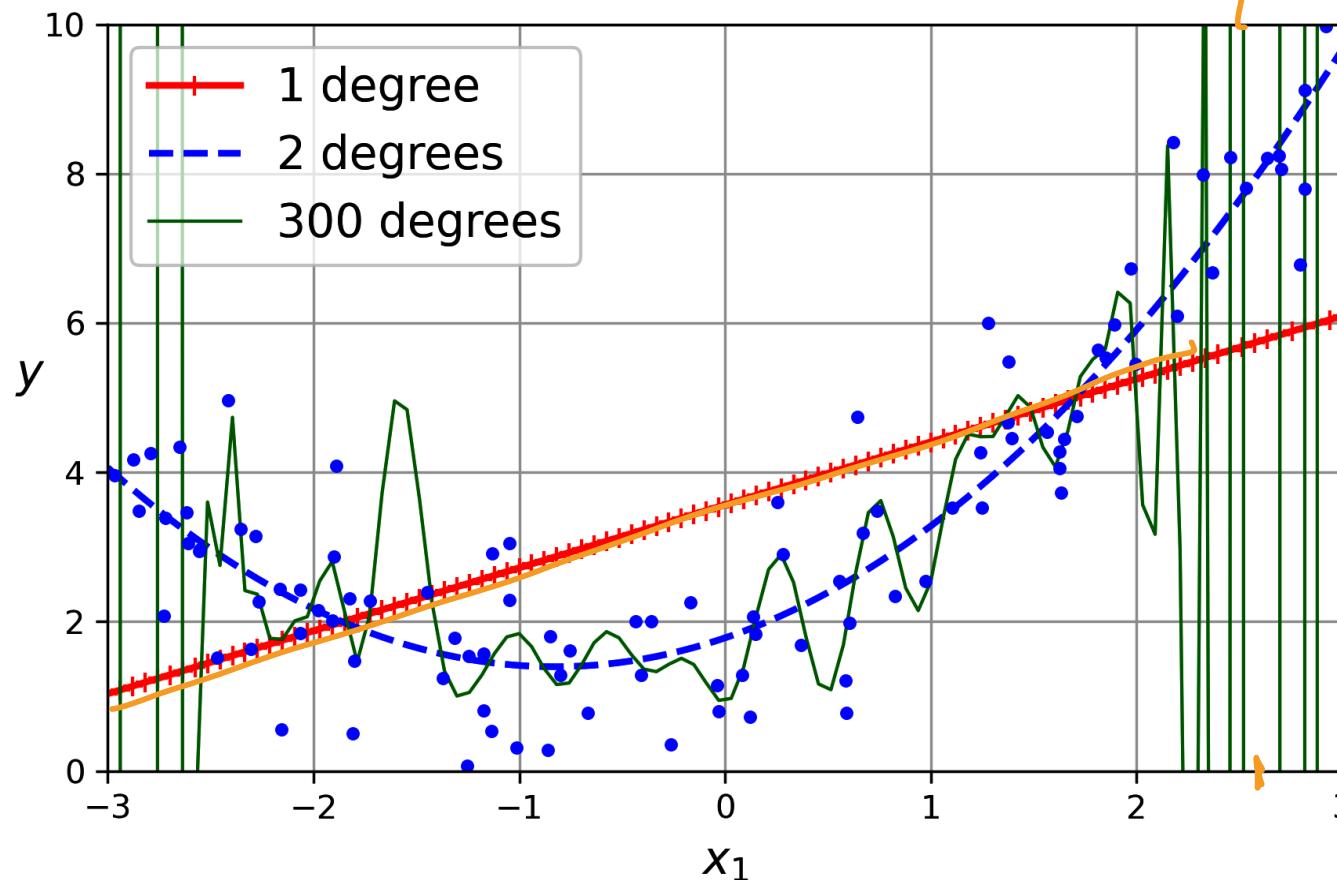
```
>>> from sklearn.preprocessing import PolynomialFeatures  
>>> poly_features = PolynomialFeatures(degree=2, include_bias=False)  
>>> X_poly = poly_features.fit_transform(X)  
>>> X[0]  
array([-0.75275929])  
>>> X_poly[0]  
array([-0.75275929,  0.56664654])  
  
>>> lin_reg = LinearRegression()  
>>> lin_reg.fit(X_poly, y)  
>>> lin_reg.intercept_, lin_reg.coef_  
(array([1.78134581]), array([[0.93366893, 0.56456263]]))
```

$$\hat{y} \quad \hat{\theta} \quad \text{estimated}$$

Not bad: the model estimates  $\hat{y} = 0.56x_1^2 + 0.93x_1 + 1.78$  when in fact the original function was  $y = 0.5x_1^2 + 1.0x_1 + 2.0 + \text{Gaussian noise.}$



# High Order Polynomial – Nice fit?

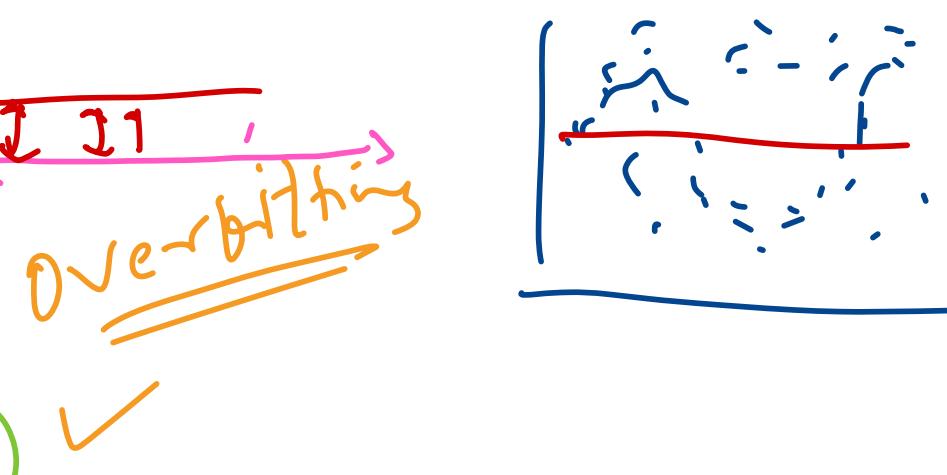
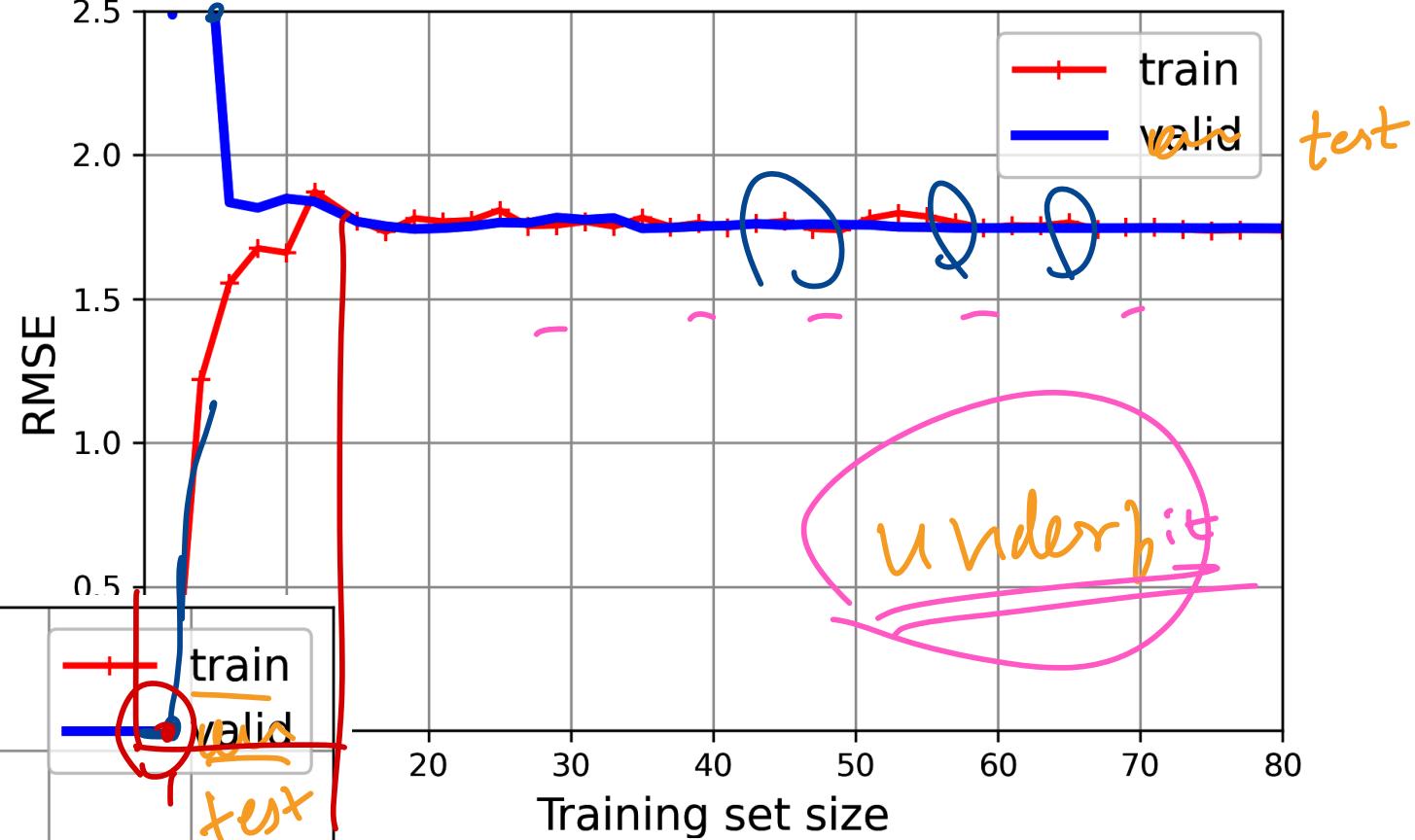
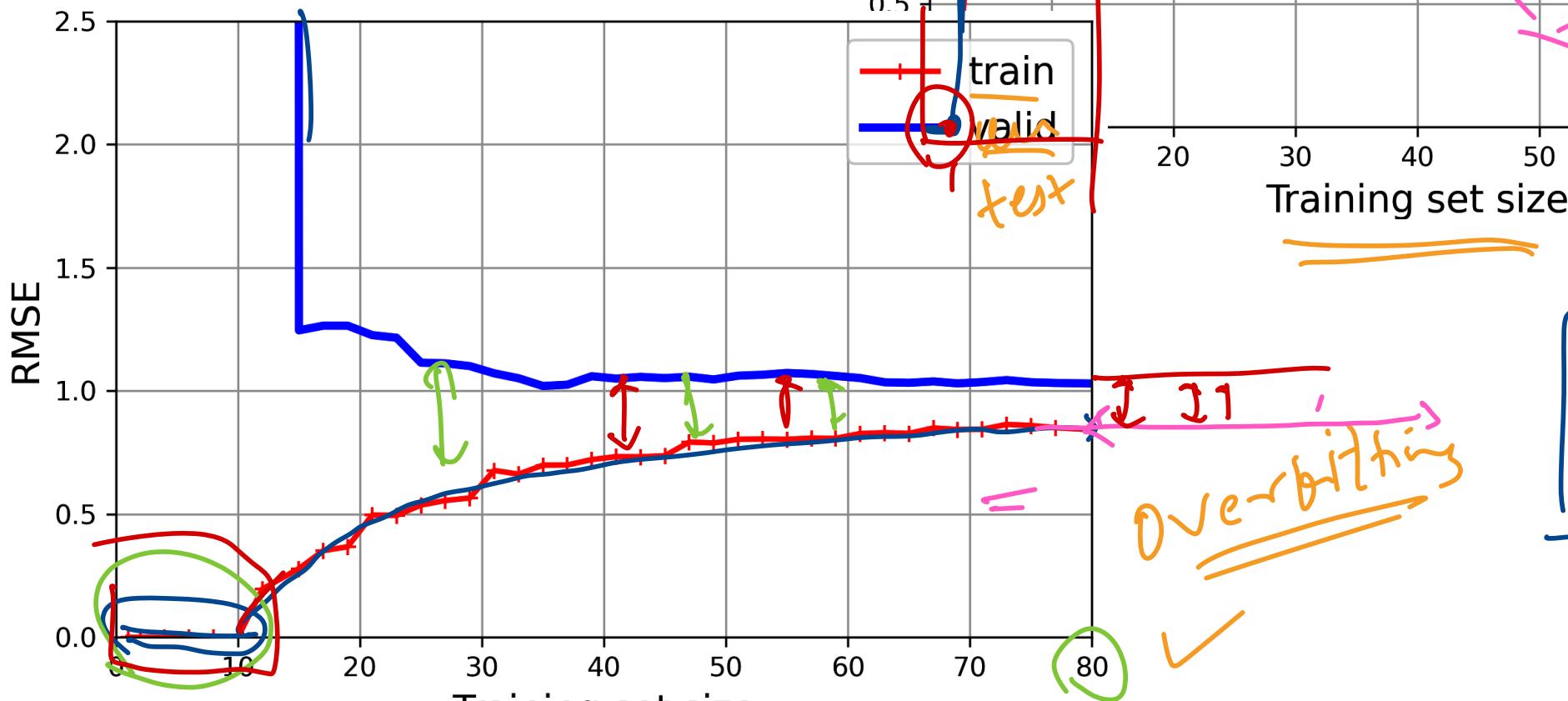


$\epsilon_{\text{var}}$   
v · low

# Learning Curves

Hands on Machine  
with Scikit

Pg 153  
154

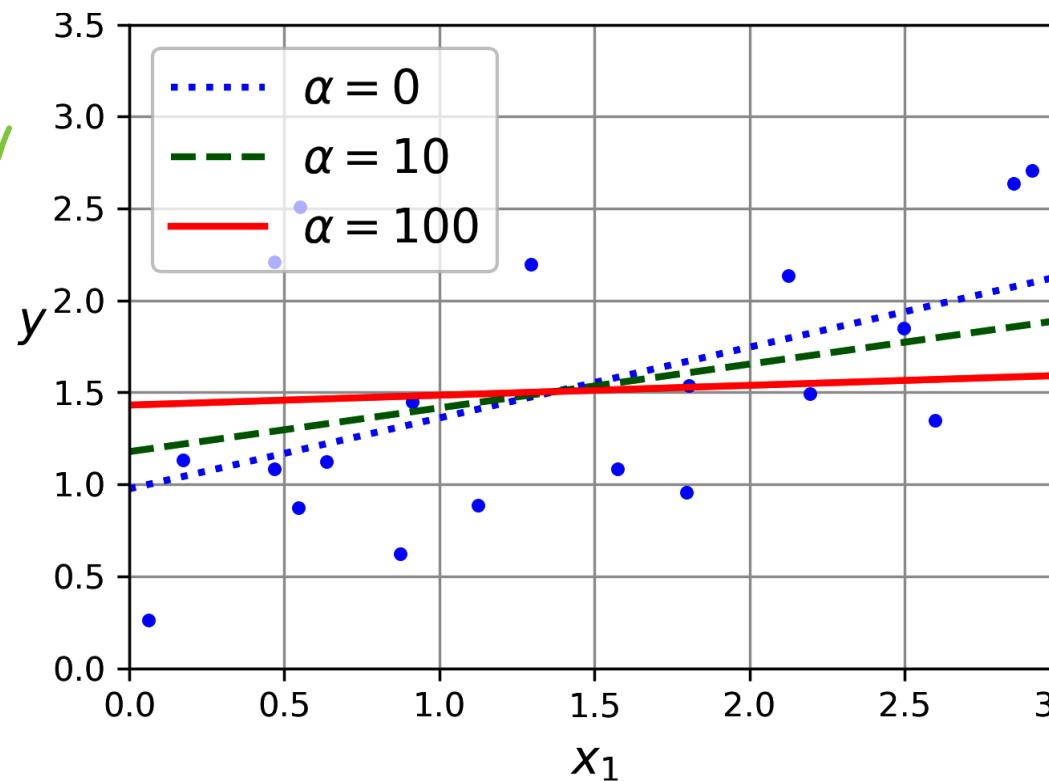


# Regularization

## Regularized Linear Models – tackles overfitting

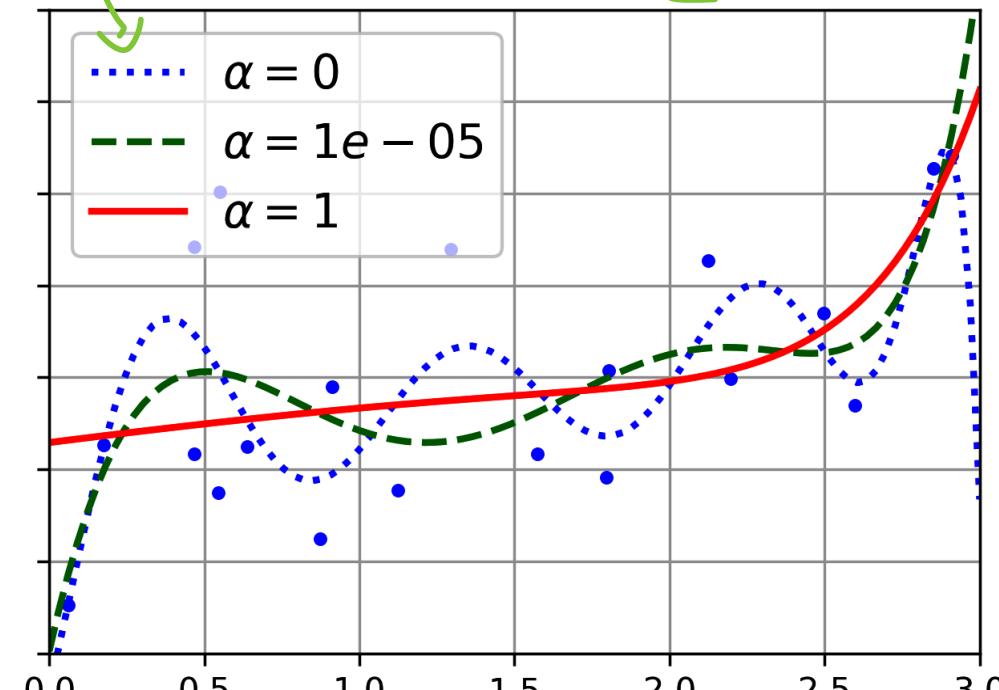
- Ridge regression

Linear and polynomial modes.



$$J(\theta) = \text{MSE}(\theta) + \frac{\alpha}{m} \sum_{i=1}^n \theta_i^2$$

$\ell_2$



- Closed form solution for Ridge Regression

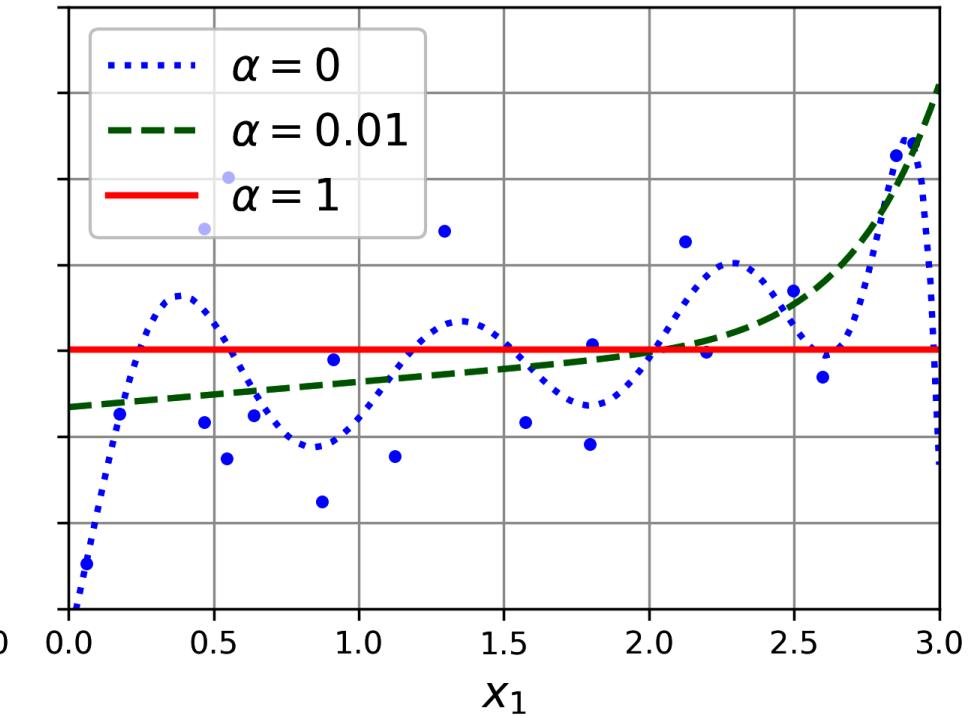
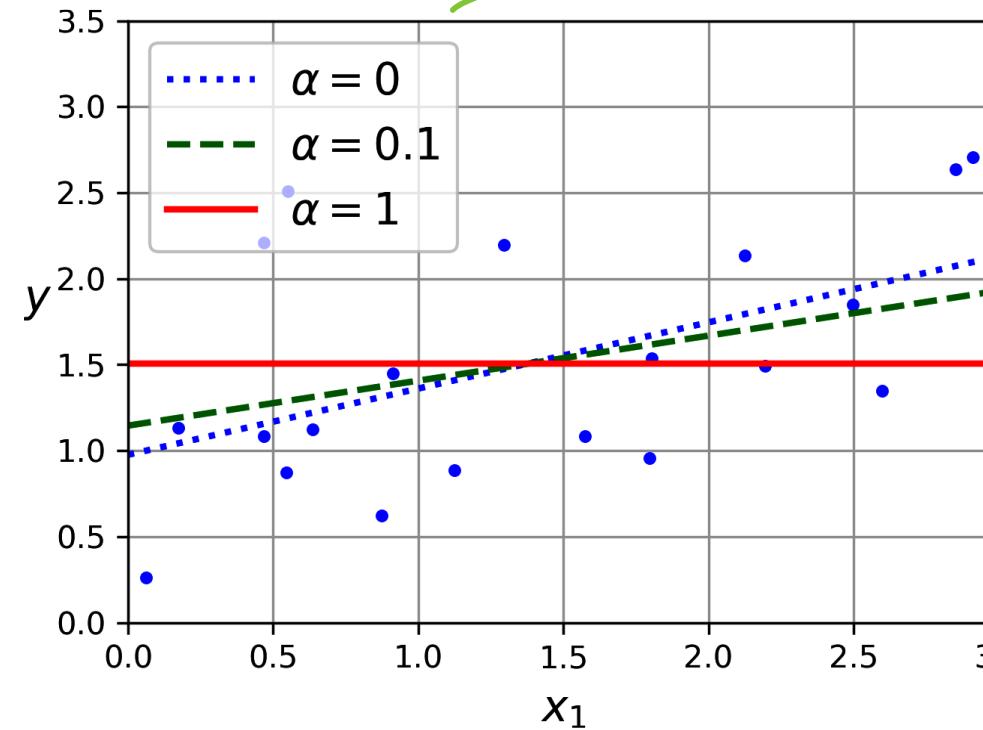
$$\hat{\theta} = (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{A})^{-1} \mathbf{X}^\top \mathbf{y}$$

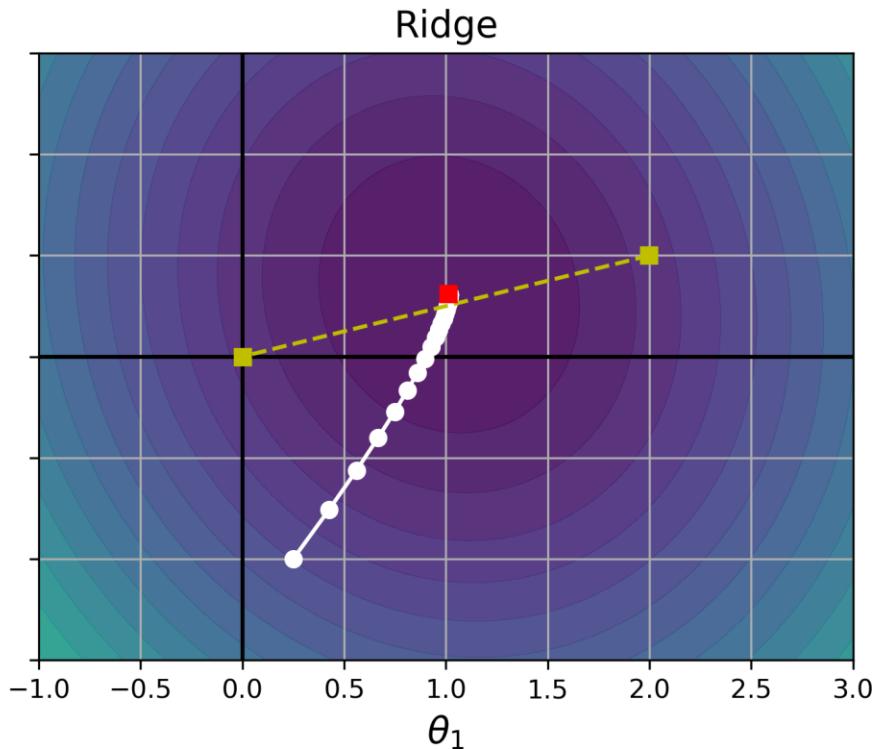
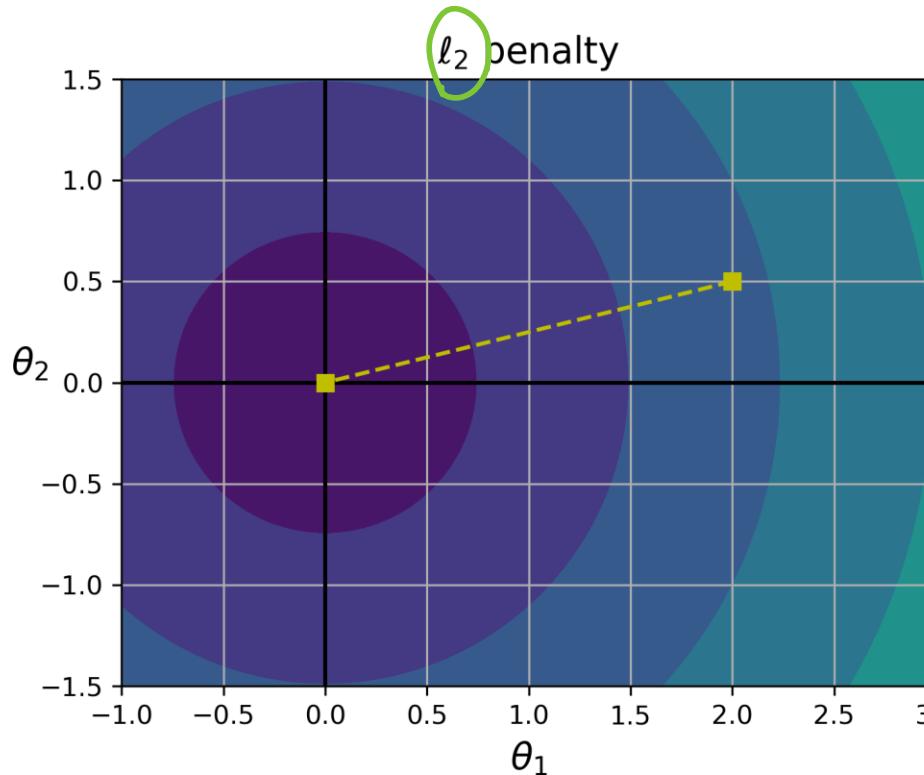
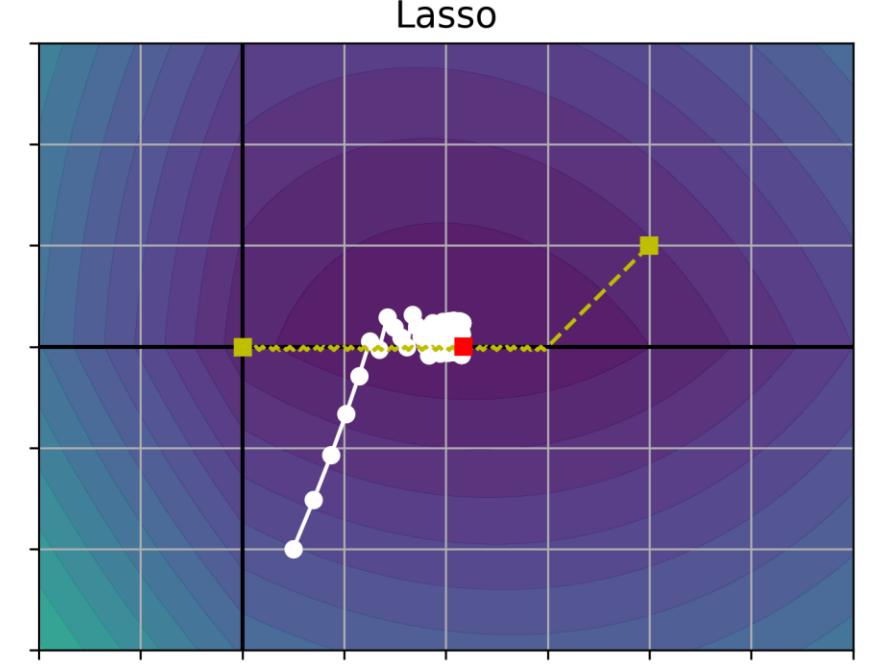
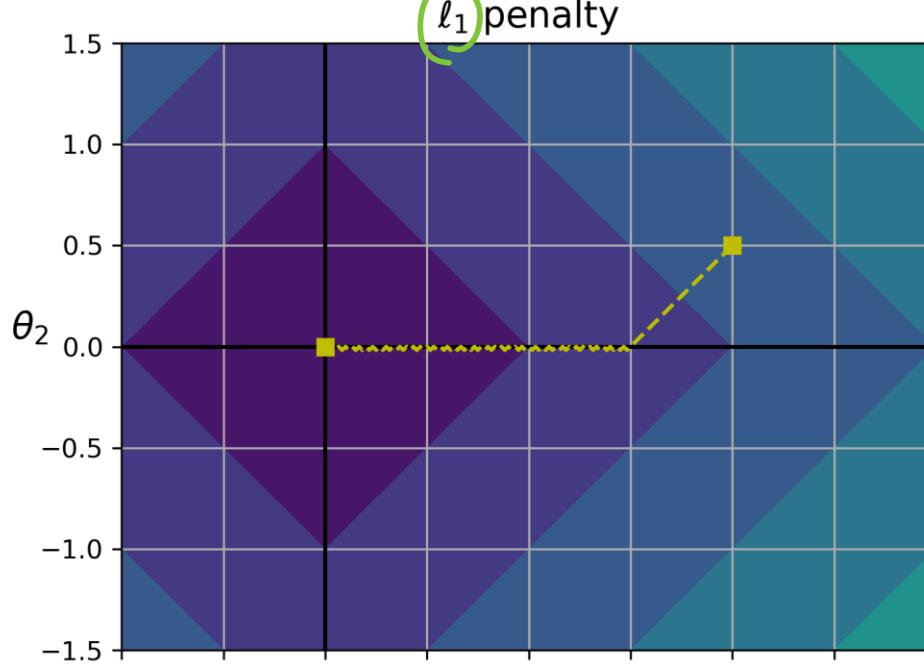
# Lasso Regression

156 - 160

$$J(\theta) = \text{MSE}(\theta) + 2\alpha \sum_{i=1}^n |\theta_i|$$

$\ell_1$





# Elastic Net Regression

- A middle ground between ridge and lasso regression
- Weighted sum of both Ridge and Lasso regression - mix ratio  $r$
- $r=0$  ... Ridge regression
- $r=1$  ... Lasso regression

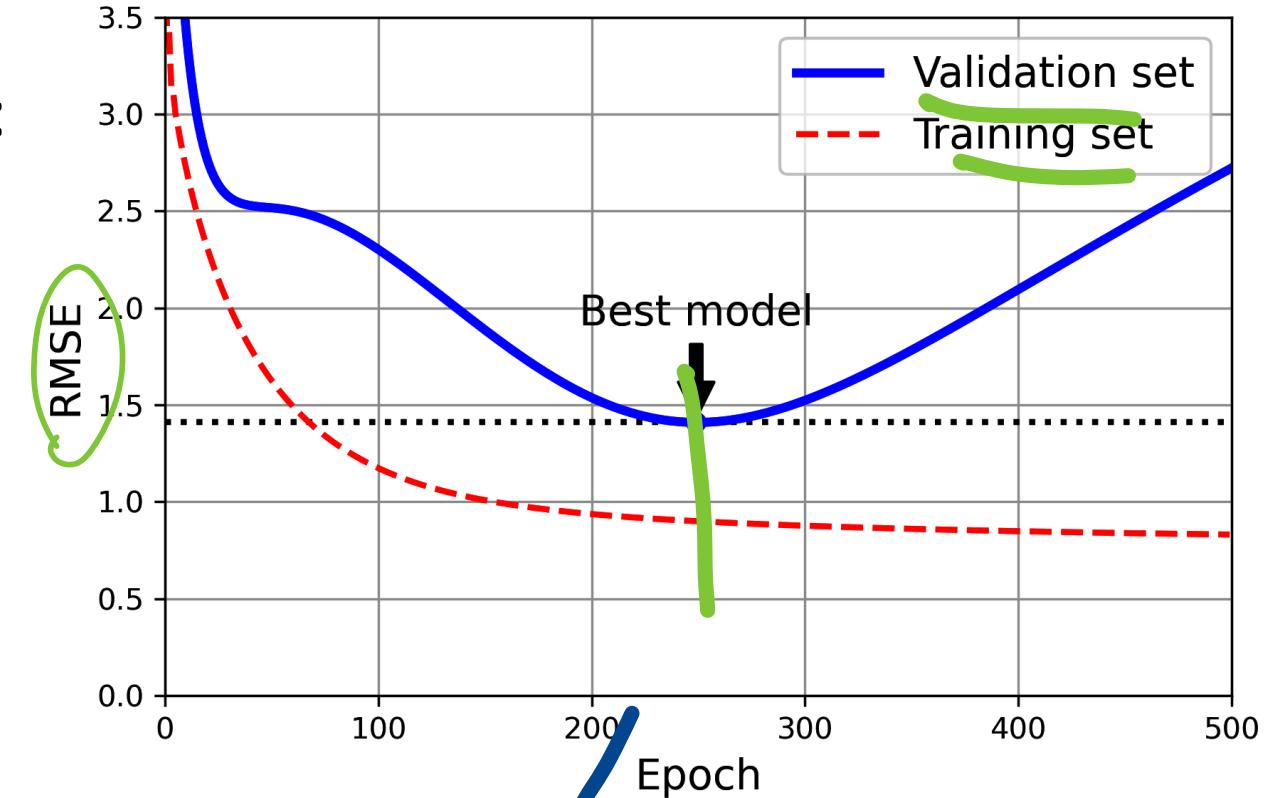
$$J(\boldsymbol{\theta}) = \text{MSE}(\boldsymbol{\theta}) + r(2\alpha \sum_{i=1}^n |\theta_i|) + (1 - r)\left(\frac{\alpha}{m} \sum_{i=1}^n \theta_i^2\right)$$

# Which Regression

- Plain, Ridge, Lasso, Elastic Net ?
- Preferable to have some regularization
- Ridge is a good default
- If you feel only a few features are useful – prefer Lasso or Elsatic Net as they reduce the weights of useless features to 0.
- Elastic net is preferred as compared to Lasso ...

# Early stopping

- Another different regularization:
- Stop training when the validation error reaches a minimum – Early stopping
- After min validation error – overfitting has started.



# Some Code ...

```
from copy import deepcopy
from sklearn.metrics import mean_squared_error
from sklearn.preprocessing import StandardScaler

X_train, y_train, X_valid, y_valid = [...] # split the quadratic dataset

preprocessing = make_pipeline(PolynomialFeatures(degree=90, include_bias=False),
                             StandardScaler())
X_train_prep = preprocessing.fit_transform(X_train)
X_valid_prep = preprocessing.transform(X_valid)
sgd_reg = SGDRegressor(penalty=None, eta0=0.002, random_state=42)
n_epochs = 500
best_valid_rmse = float('inf')

for epoch in range(n_epochs):
    sgd_reg.partial_fit(X_train_prep, y_train)
    y_valid_predict = sgd_reg.predict(X_valid_prep)
    val_error = mean_squared_error(y_valid, y_valid_predict, squared=False)
    if val_error < best_valid_rmse:
        best_valid_rmse = val_error
        best_model = deepcopy(sgd_reg)
```