

Unsupervised Learning: unlabeled data

Clustering algorithm: DNA Microarray

grouping customer
similar Data point together

Dimensional Reduction: compress data to fewer numbers

Anomaly detection: Find unusual Data point

Terminology

Training set: Data used to train the model

$x \rightarrow$ input variable
 $y \rightarrow$ output variable
(target variable)

Training set → features
target

Learning algorithm

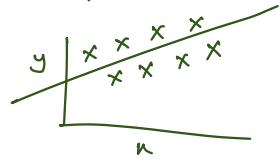
$n \rightarrow f \rightarrow g$
function

feature model prediction
size → f → price

How to represent?

$$f_{w,b}(x) = wx + b$$

$$f(x) =$$



$$f(w,b)(x) = wx + b$$

$$f(x) = wx + b$$

linear regression with one variable

Cost Function:

w, b do → parameter

$$\sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 \text{ error} = J(w, b) \rightarrow \text{squared error cost function}$$

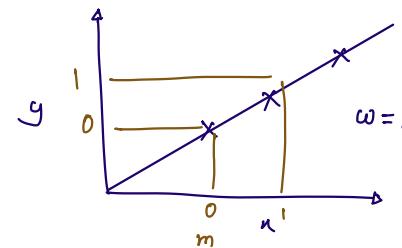
m = number of training examples

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m f(w, b)(x^{(i)}) - y^{(i))^2}$$

$w, b \rightarrow$ can be adjusted

Cost Function Intuition:

$f_w(x)$
For fixed w , function of x



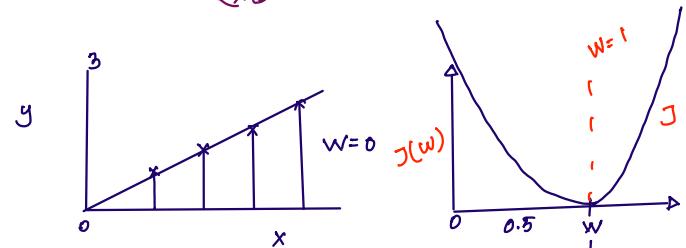
$$J(1) = \frac{1}{2m} \sum_{i=1}^m f(1)(x_i - g_i)^2$$

$$= \frac{1}{2m} (0+0)$$

$$= 0$$

$$J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

$$= \frac{1}{2 \times 3} [3.5] = 0.58$$

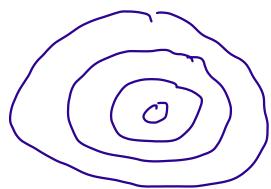


Goal → minimize $J(w)$

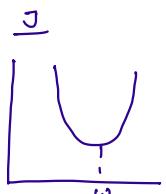
general case → minimize $J(w, b)$

f_w, b

cost function



$J(w, b)$



b

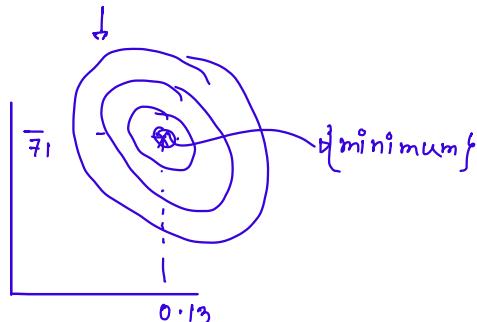
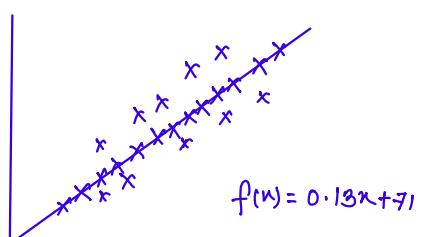
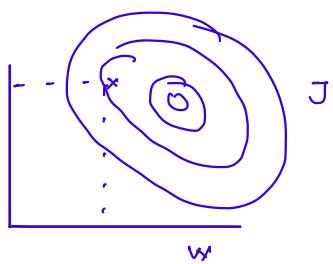
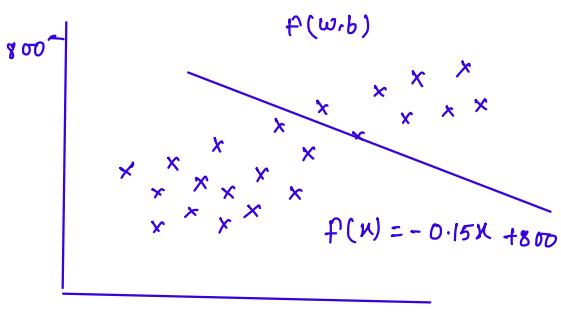


$J(w=10, b=15)$



④ contour plot

④



- Best fit line → Manually not possible

Cost function,

$$J(w, b) = \frac{1}{2m} \sum_{i=0}^{m-1} f(w, b)(x^{(i)} - y^{(i)})^2$$

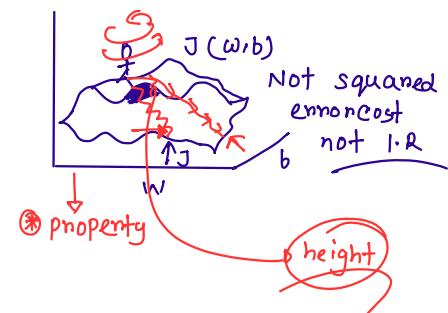
$$f(w, b)(x^{(i)}) = w x^{(i)} + b$$

Gradient descent

↳ Not bound in
for linear regression
 \downarrow
 $\min J(w, b) \rightarrow$ To minimize
 w, b

ML also
in DL

start with some w, b
 \downarrow set $(w=0, b=0)$
keep changing w, b to
reduce $J(w, b)$

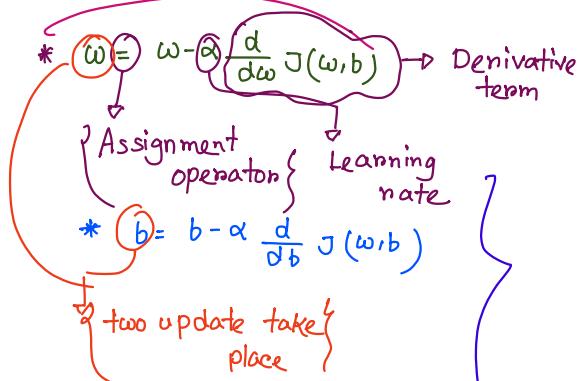


Training L.R.

Implementing Gradient

Descent

Gradient Descent Algorithm:



④ simultaneously update w, b

$$\begin{aligned} \text{temp_}w &= w - \alpha \frac{\partial}{\partial w} J(w, b) && \text{Not partial} \\ \text{temp_}b &= b - \alpha \frac{\partial}{\partial b} J(w, b) && \text{Hence full} \\ w &= \text{temp_}w && \text{Mathematical} \\ b &= \text{temp_}b \end{aligned}$$

Regression with multiple input variable

$$f_{\text{wrb}}(x) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

$$\bar{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n]$$

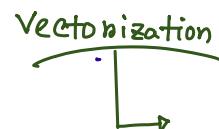
b is a number

$$\bar{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$$

multiple linear regression

$$f_{\bar{w}, b}(\bar{x}) = \bar{w} \cdot \bar{x} + b = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

(dot product)



→ Makes code faster!

$$f_{\bar{w}, b}(\bar{x}^d) = \sum_{j=1}^n w_j x_j + b$$

#vectorization:

$$f = \text{np.dot}(w^T x) + b \quad \left\{ \begin{array}{l} j = 1 \dots n \\ 1, 2, 3 \end{array} \right\}$$

if $f = 0$
 for j in range(0, n):
 $f = f + w[j] * x[j]$
 $f = f + b$

Without vectorization

Gradient Descent:

$$w_1 = w_1 - \alpha \cdot d_1$$

$$w_2 = w_2 - \alpha \cdot d_2$$

$$\vdots \quad \vdots$$

$$w_n = w_n - \alpha \cdot d_n$$

{Without vectorization}

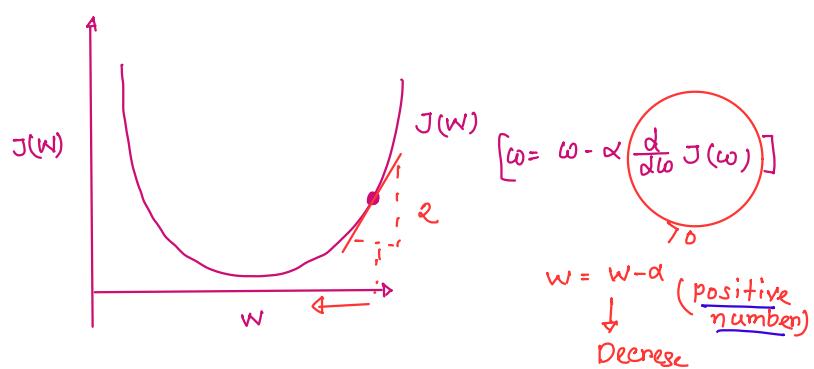
$$\text{for } j \text{ in range}(0, n)$$

$$w[j] = w[j] - \alpha \cdot d[j]$$

$$[w = w - \alpha \cdot d] \rightarrow \text{With vectorization}$$

Normal equation → alternate
 ↳ For linear regression method

gradient descent → recommended Method



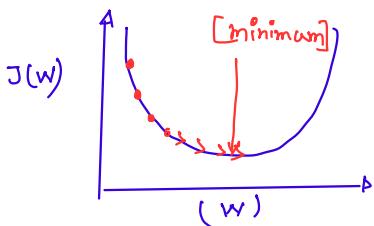
$$\frac{d}{dw} J(w) < 0$$

$$w = w - \alpha \quad \downarrow \quad \text{(negative number)}$$

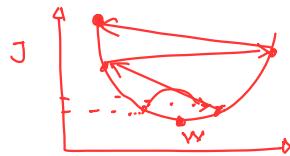
Increase

Learning Rate

if α is too small...

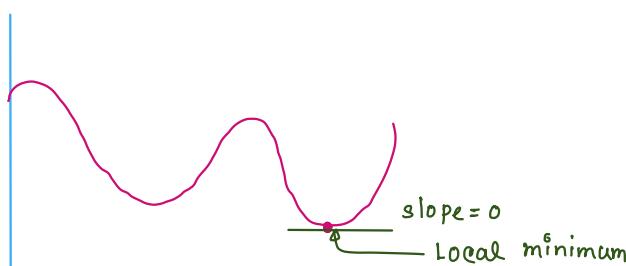


too Large:

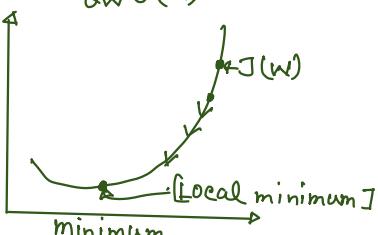


Gradient descent may:

- Overshoot
- Never reach minimum



$$w = w - \alpha \frac{d}{dw} J(w)$$

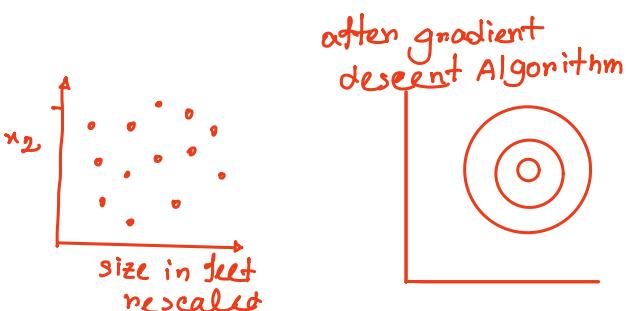
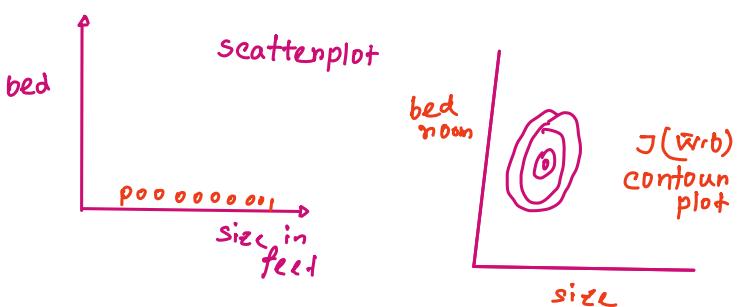


Feature and Parameter value

$$\text{price} = w_1 x_1 + w_2 x_2 + b$$

↓ ↓
size bedrooms

Features



Feature scaling

$$3000 \leq x_1 \leq 2000$$

$$x_1 \text{ scaled} = \frac{x_1 - \mu_1}{\sigma_1}$$

make range small

Mean Normalization:

$$\mu_1 = \frac{x_1 - \mu_1}{2000 - 300}$$

$$-0.18 \leq x_1 \leq 0.82$$

purpose: center the Data
Z score Normalization → standard score

standard deviation $\sigma_1 = 950$ $\sigma_2 = 1.9$

$$x_1 = \frac{x_1 - \mu_1}{\sigma_1}$$

values centered with std. dev. (1)

$$-0.67 \leq x_1 \leq 3.1$$

Feature scaling

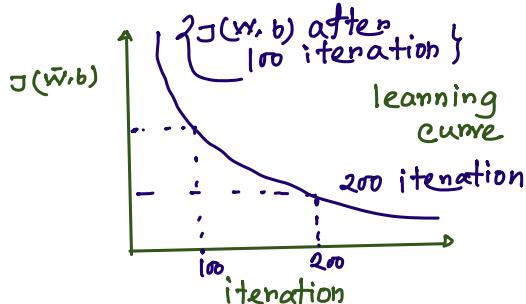
aim for about $-1 \leq x_j \leq 1$ for each feature

$-3 \leq x_j \leq 3$

$-0.3 \leq x_j \leq 0.3$

Make sure gradient descent is working correctly.

$$\text{objective : } \min J(\bar{w}, b)$$



$J(\bar{w}, b)$ should decrease after every iteration

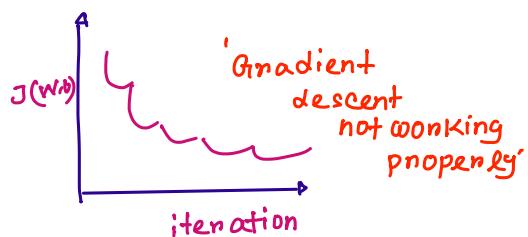
Automatic convergence test:

Let ϵ be $10^{-3} 0.001$

$J(\bar{w}, b)$ decreases by $\leq \epsilon$ in one iteration.

declare convergence

"Choose Learning Rate"

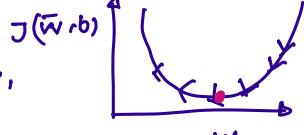


Adjusting Learning Rate

choosing Learning Rate!



smaller Learning rate

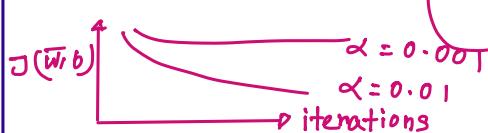


$$w_1 = w_1 + \alpha d_1 \quad X$$

$$w_1 = w_1 - \alpha d_1 \quad \checkmark$$

$J(\bar{w}, b)$, α should decrease on every iteration

values of α try to $0.001, 0.01, 0.1, 1 \dots$



Feature Engineering

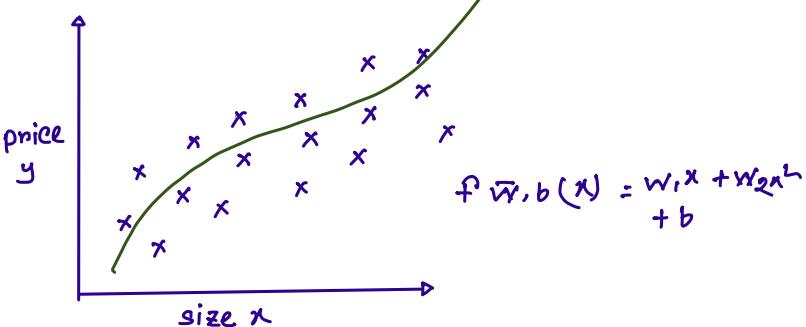
$$f(\bar{w}, b)(x) = \underline{w_1 x_1} + \underline{w_2 x_2} + b$$

area = frontage \times depth

$$x_3 = x_1 x_2 \rightarrow \text{new feature}$$

transforming,
combining
original
features.

Polynomial Regression



Classification

yes/no

— False/true

Binary Classification

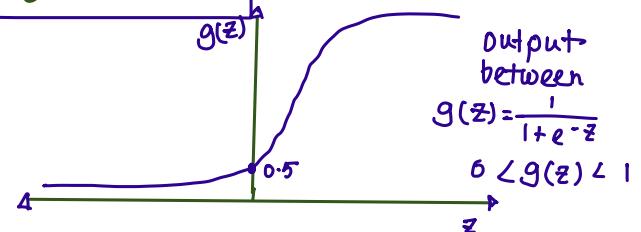


[Class - category]

- decision boundary

Logistic Regression

Sigmoid Function:



$$\begin{aligned} f(\bar{w}, b)(\bar{x}) &= g(\underbrace{\bar{w} \cdot \bar{x} + b}_z) \\ &= \frac{1}{1+e^{-(\bar{w} \cdot \bar{x} + b)}} \end{aligned}$$

Probability:

$$[\# p(y=0) + p(y=1) = 1]$$

Decision Boundary

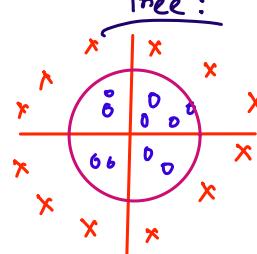
$$P(y=1|\bar{x}; \bar{w}, b)$$



Decision Boundary

$$\begin{aligned} z &= \bar{w} \cdot \bar{x} + b = 0 \\ z &= x_1 + x_2 - 3 = 0 \end{aligned}$$

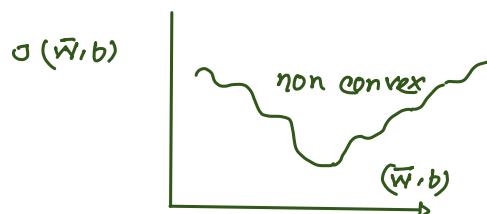
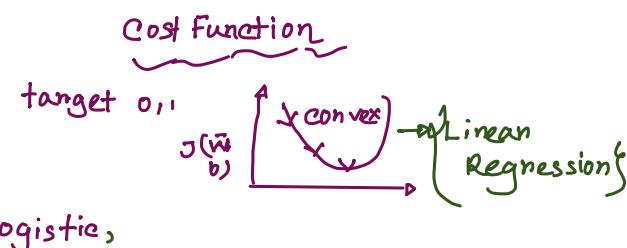
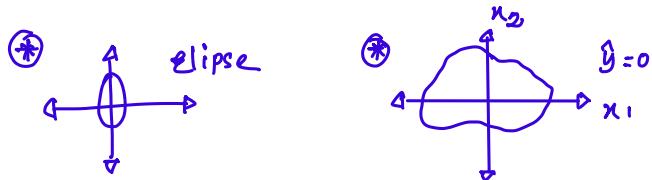
Non Linear Decision Tree:



$$z = x_1 + x_2 - 1 = 0$$

$$f(\bar{w}, b)(\bar{x}) = g(z)$$

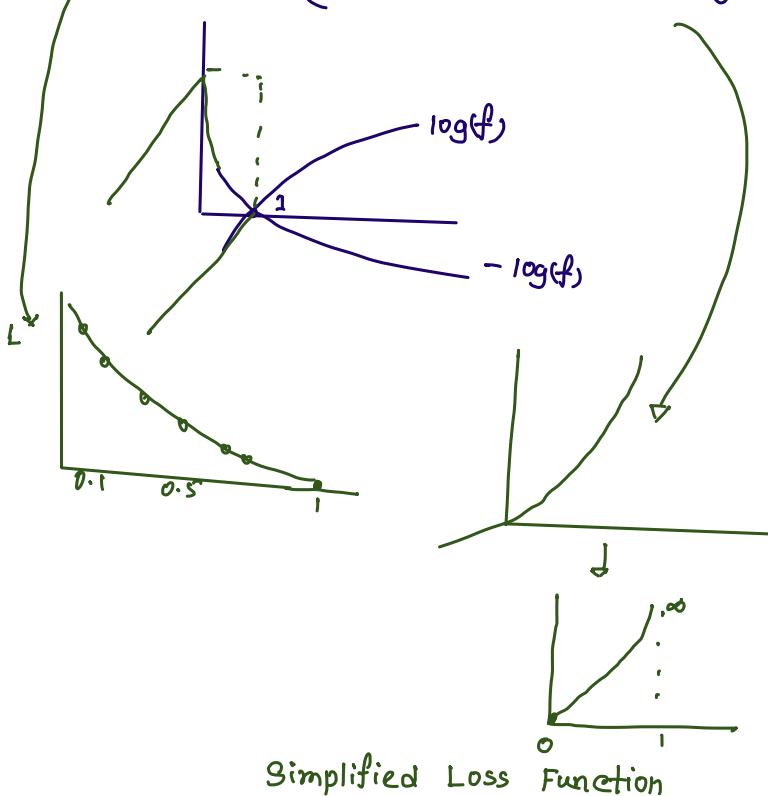
$$\begin{aligned} x_1 + x_2 &= 1 \\ \text{circle} &= g(w_1 x_1 + w_2 x_2) \end{aligned}$$



logistic Loss function:

$$L f(\bar{w}, b(\bar{x}(i), y(i)) = \begin{cases} -\log f(\bar{w}, b(\bar{x}(i))) & \text{if } y(i) = 1 \\ -\log (1 - f(\bar{w}, b(\bar{x}(i))) & \text{if } y(i) = 0 \end{cases}$$

↑ can reach global convex min.



$$L f(\bar{w}, b(\bar{x}(i), y(i)) = -y(i) \log f(\bar{w}, b(\bar{x}(i))) + (1-y(i)) \log (1 - f(\bar{w}, b(\bar{x}(i))))$$

if, $y(i) = 1$

$L f(\bar{w}, b(\bar{x}(i), y(i)) = -\log f(\bar{x})$

if, $y(i) = 0$

$L f(\bar{w}, b(\bar{x}(i), y(i)) = \dots$

Simplified cost Function

$$\begin{aligned} J(\bar{w}, b) &= \frac{1}{m} \sum_{i=1}^m [L f(\bar{w}, b(\bar{x}(i), y(i))] \\ &= -\frac{1}{m} \sum_{i=1}^m y(i) \log f(\bar{w}, b(\bar{x}(i))) \\ &\quad + (1-y) \log (1 - f(\bar{w}, b(\bar{x}(i)))) \end{aligned}$$

1 iteration - loss \rightarrow lennon
average loss - cost
lennon

Gradient descent

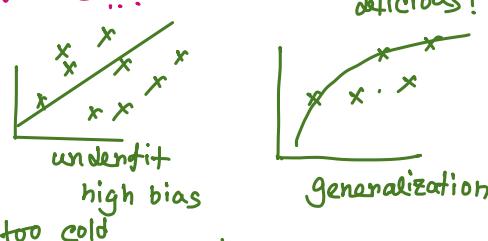
- Cost minimum

Formula same as linear regression

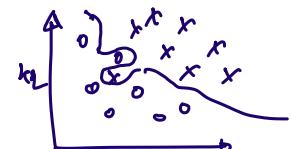
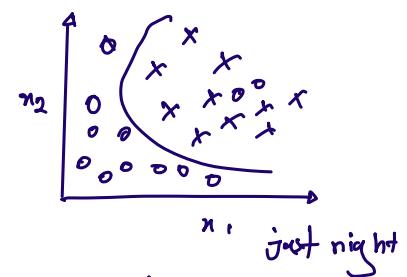
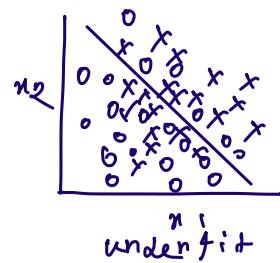
↓ significance different

** Overfitting !!!

Reg.



Clas.



Regularization:

\downarrow $\lambda \rightarrow$ Regularized term

$$\lambda > 0$$

$$\min J(\bar{w}, b) = \min_{\bar{w}, b} \sum_{i=1}^m f(\bar{w}, b(\bar{x}(i) - y(i))) + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

mean squared error

