

Supervised Machine Learning: Regression and Classification

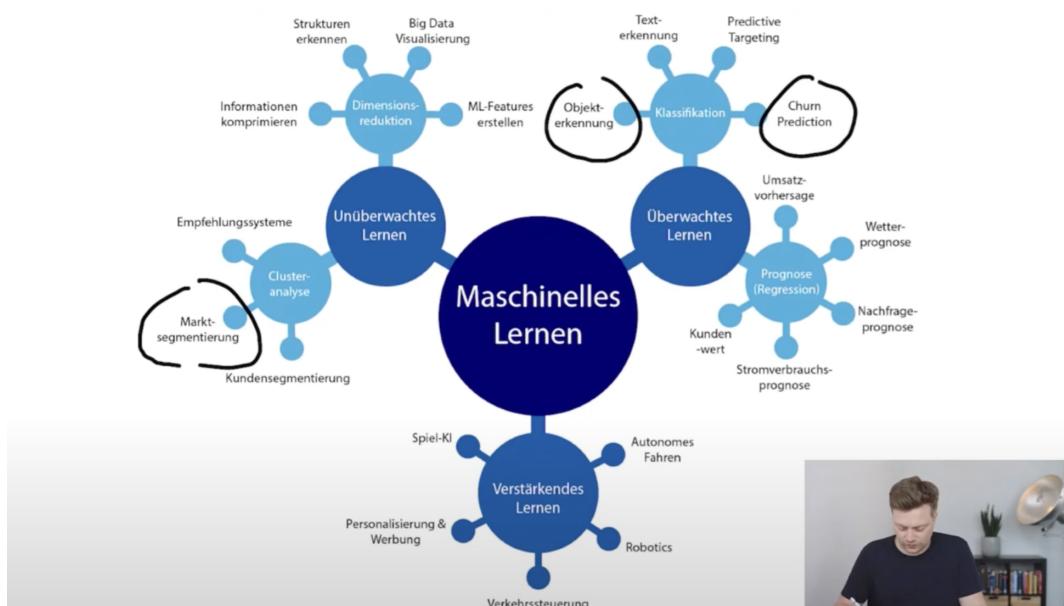
Course name: Supervised Machine Learning: Regression and Classification

Institution: Deeplearning.ai

By: Vadhna Samedy Hun

I. Introduction (Grundlage)

Machine Learning is the field of study that gives computers the ability to learn without being explicitly programmed. (Maschinelles Lernen is das Studienfach, das Computern die Fähigkeit gibt, zu lernen, ohne explizit programmiert zu werden.)



II. Machine Learning

2.1. Supervised Learning

Supervised learning



Learns from being given "right answers"

Input (X)	Output (Y)	Application
email	→ spam? (0/1)	spam filtering
audio	→ text transcripts	speech recognition
English	→ Spanish	machine translation
ad, user info	→ click? (0/1)	online advertising
image, radar info	→ position of other cars	self-driving car
image of phone	→ defect? (0/1)	visual inspection

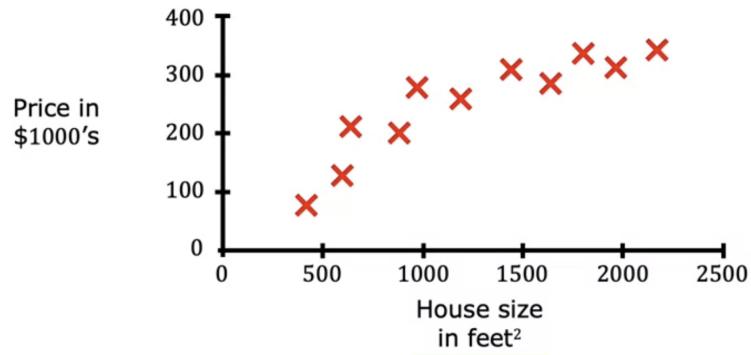
- ▼ There are 2 types of supervised learning:

Regression: predict continuous value

Classification (Logistic Regression): predict categories (two categories)

Learn from data labeled with the "right answers".

Regression: Housing price prediction

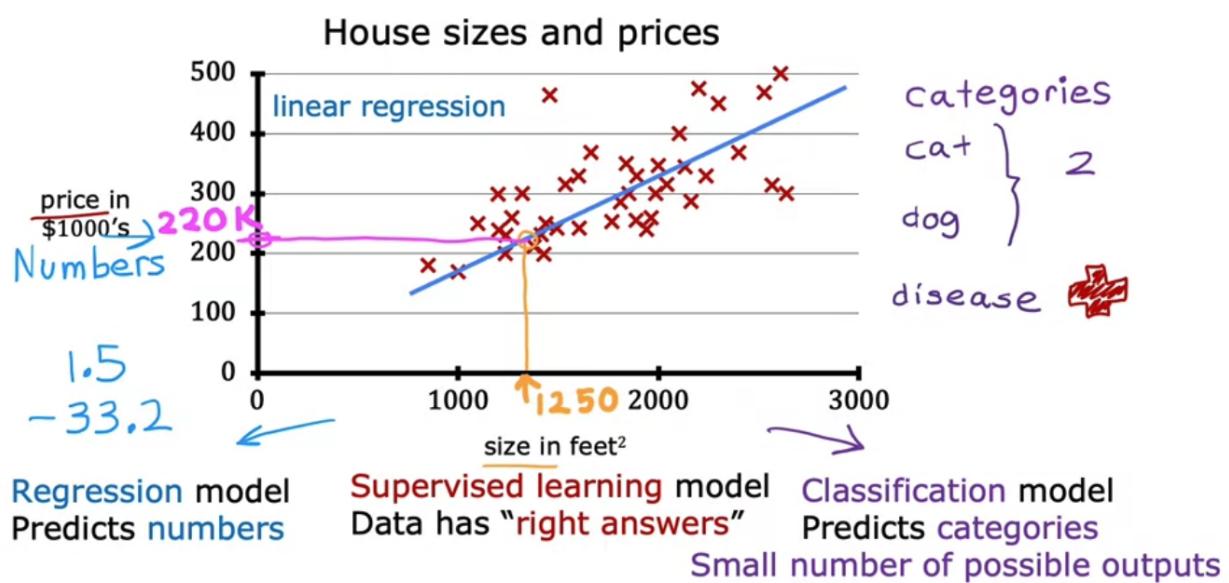


Stanford ONLINE

DeepLearning.AI

Andrew Ng

2.1.1. Regression



Stanford ONLINE

DeepLearning.AI

Andrew Ng

Terminology

Training set:	x size in feet ²	y price in \$1000's
	→	→
(1) 2104		400
(2) 1416		232
(3) 1534		315
(4) 852		178
...		...
(47) 3210		870
	$x = 2104$	$y = 400$
	(x, y)	$(2104, 400)$

Notation:

x = "input" variable
feature

y = "output" variable
"target" variable

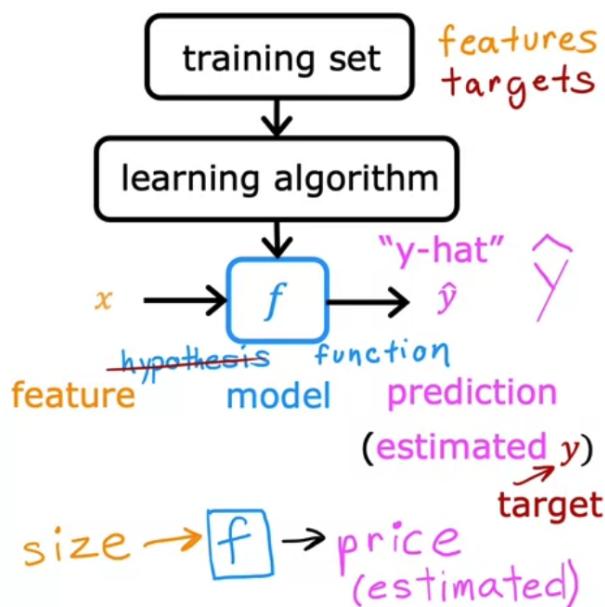
m = number of training examples

(x, y) = single training example

$(x^{(i)}, y^{(i)})$

$(x^{(i)}, y^{(i)})$ = i^{th} training example

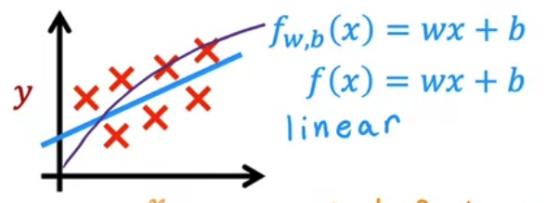
(1st, 2nd, 3rd ...)



How to represent f ?

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

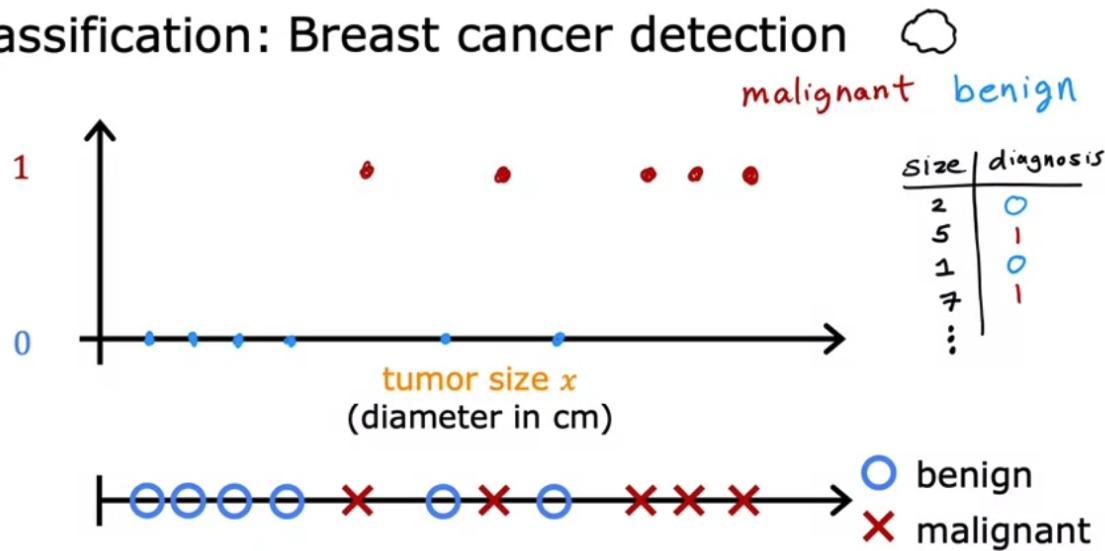
$$f(x)$$



single feature x
Linear regression with one variable.
size
one variable
Univariate linear regression.

2.1.2. Classification

Classification: Breast cancer detection



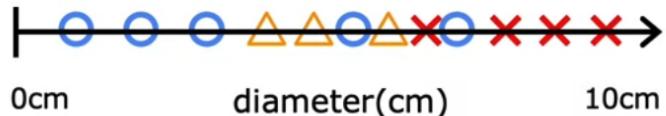
Stanford ONLINE

DeepLearning.AI

Andrew Ng

Classification: Breast cancer detection

- benign
- ✗ malignant type 1
- △ malignant type 2



class category

Classification

predict categories cat dog benign malignant 0, 1, 2

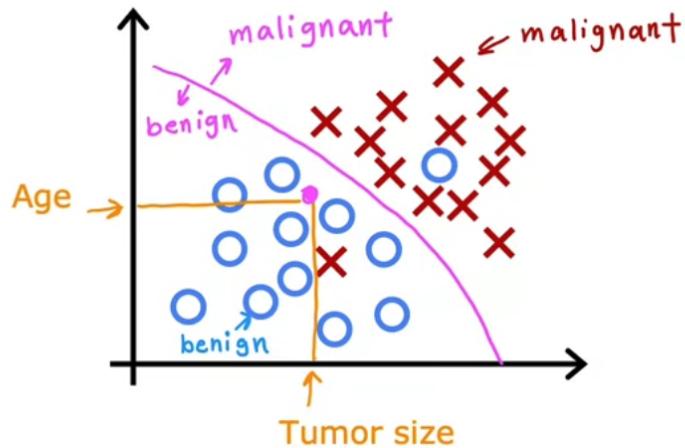
small number of possible outputs

Stanford ONLINE

DeepLearning.AI

Andrew Ng

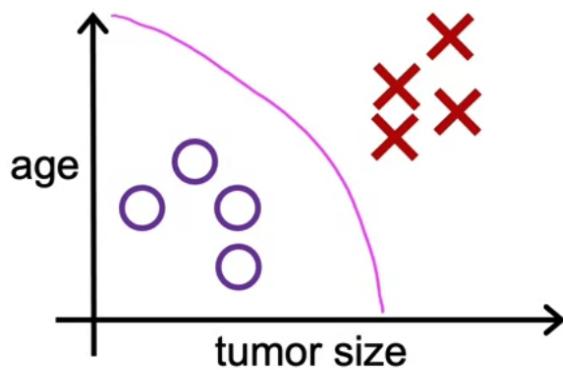
Two or more inputs



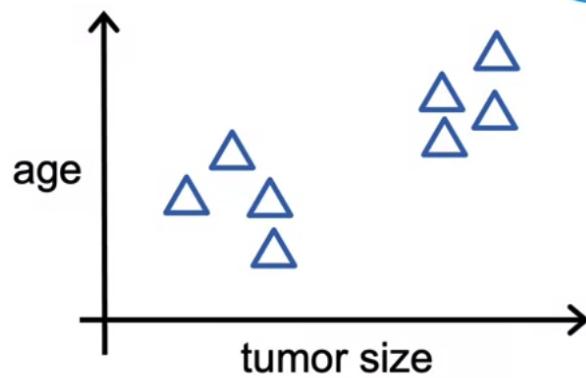
2.2. Unsupervised Learning

Find something interesting in unlabeled data.

Supervised learning
Learn from data **labeled**
with the “**right answers**”



Unsupervised learning
Find something interesting
in **unlabeled** data.

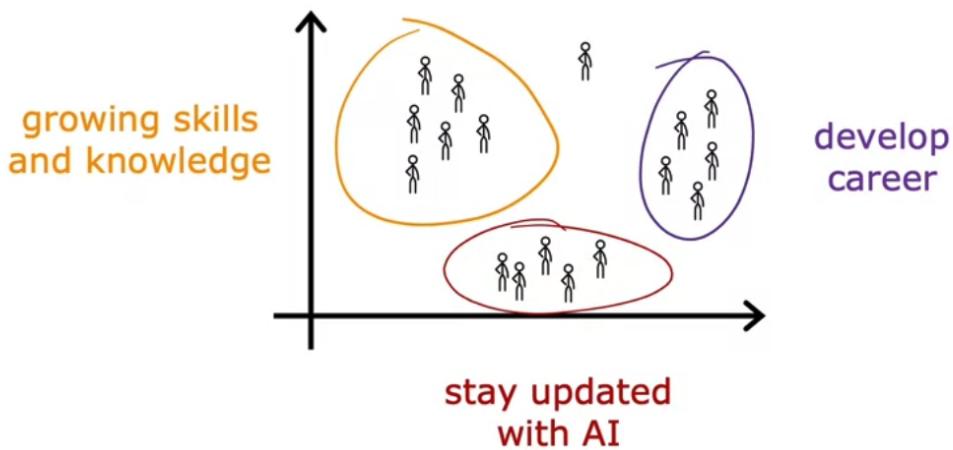


Stanford | ONLINE

DeepLearning.AI

Andrew Ng

Clustering: Grouping customers



Stanford | ONLINE

DeepLearning.AI

Andrew Ng

Data only comes with inputs x , but not output labels y . Algorithm has to find structure in the data.

▼ Clustering

Group similar data points together.

▼ Anomaly detection

Find unusual data points.

▼ Dimensionality reduction

Compress data using fewer numbers.

Unsupervised learning

Data only comes with inputs x , but not output labels y .

Algorithm has to find **structure** in the data.

Clustering

Group similar data points together.

Dimensionality reduction

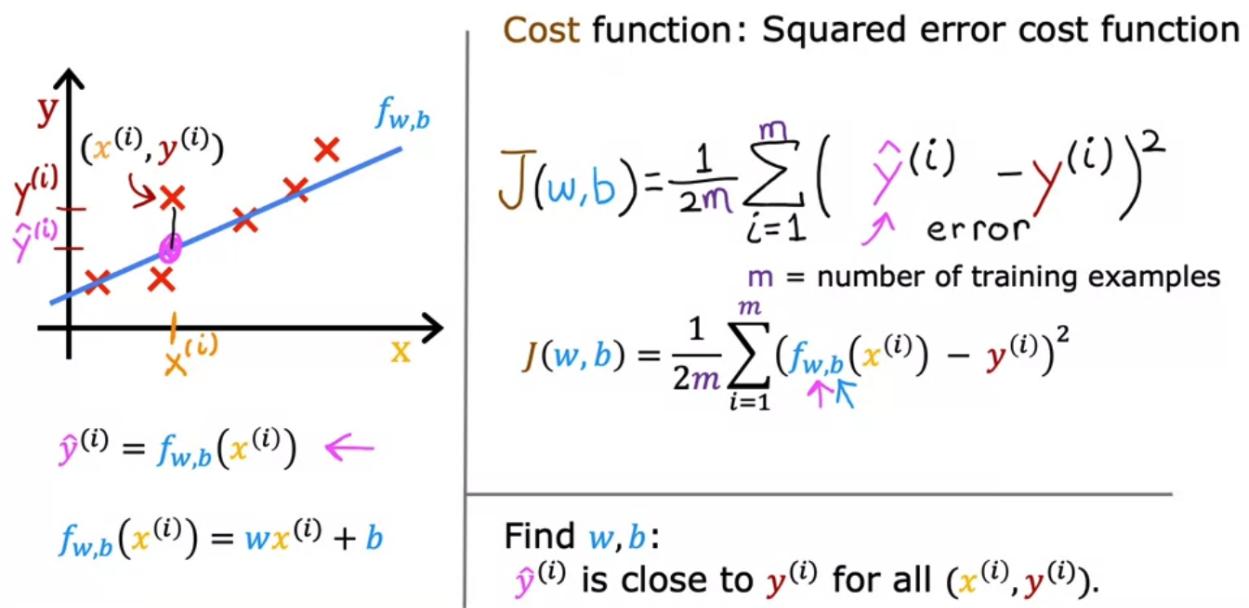
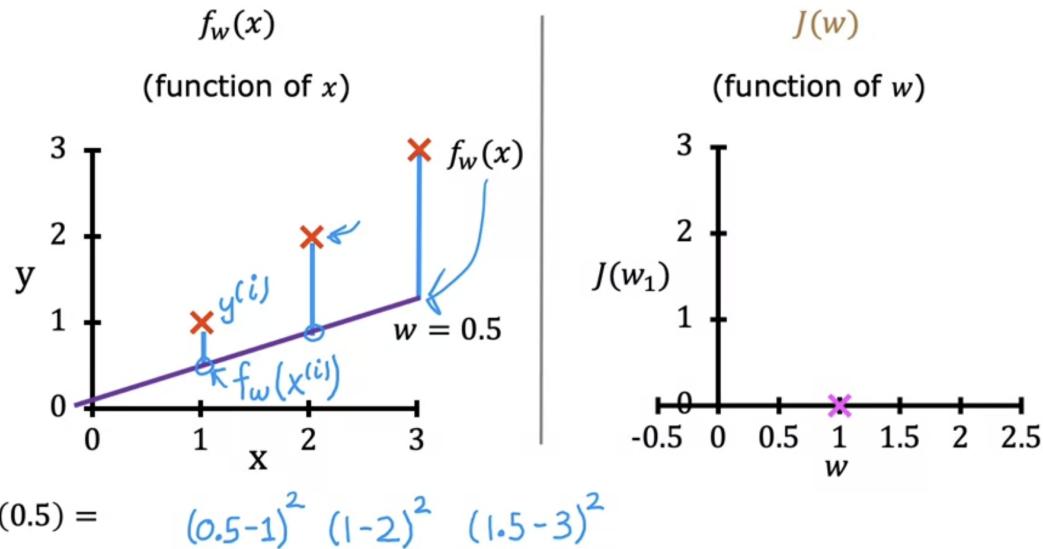
Compress data using fewer numbers.

Anomaly detection

Find unusual data points.

2.3. Cost function

Cost function is telling us how well our model is doing.



model:

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

parameters:

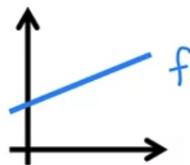
$$w, b$$

cost function:

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

goal:

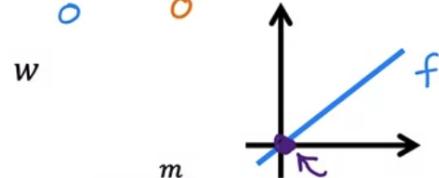
$$\underset{w, b}{\text{minimize}} J(w, b)$$



simplified

$$f_w(x) = wx$$

$$b = \emptyset$$



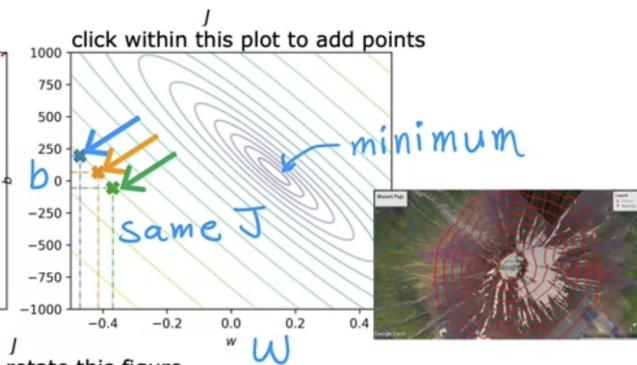
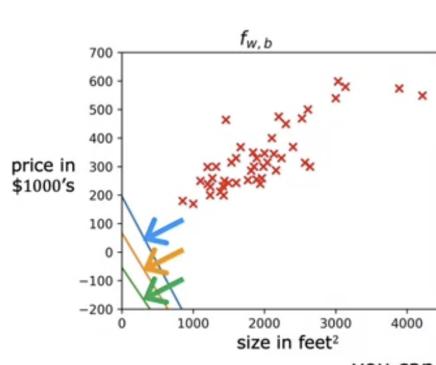
$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

$$\underset{w}{\text{minimize}} J(w)$$

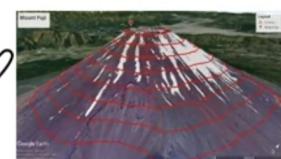
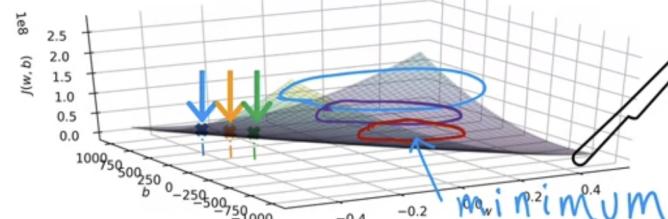
Stanford | ONLINE

DeepLearning.AI

Andrew Ng



you can rotate this figure



Stanford | ONLINE

DeepLearning.AI

Andrew Ng

Have some function $J(w, b)$ for linear regression or any function

Want $\min_{w, b} J(w, b)$

$\min_{w_1, \dots, w_n, b} J(w_1, w_2, \dots, w_n, b)$

Outline:

Start with some w, b (set $w=0, b=0$)

Keep changing w, b to reduce $J(w, b)$

Until we settle at or near a minimum

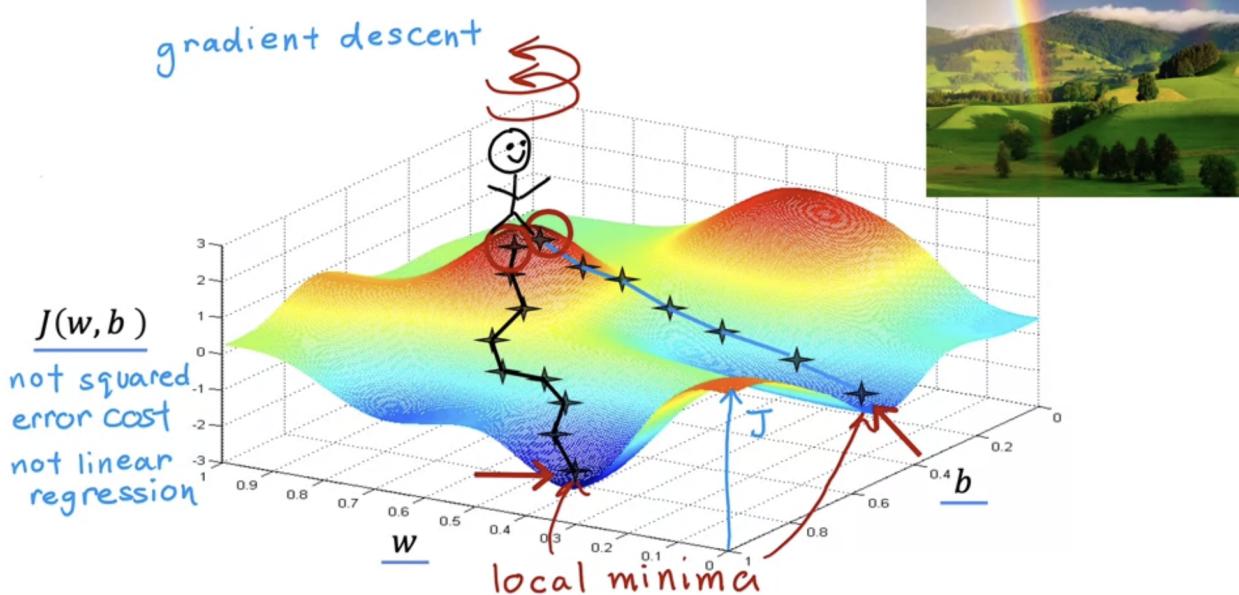
may have >1 minimum



Stanford | ONLINE

DeepLearning.AI

Andrew Ng



Stanford | ONLINE

DeepLearning.AI

Andrew Ng

SupervisedMachineLearning-WEEK-1-QUIZ-2

2.3.1. Gradient Descent

Have some function $J(w, b)$ for linear regression
or any function

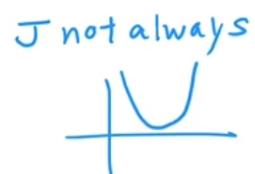
Want $\min_{w, b} J(w, b)$ $\min_{w_1, \dots, w_n, b} J(w_1, w_2, \dots, w_n, b)$

Outline:

Start with some w, b (set $w=0, b=0$)

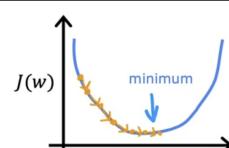
Keep changing w, b to reduce $J(w, b)$

Until we settle at or near a minimum

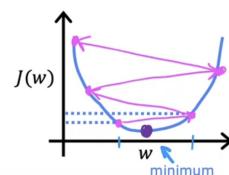


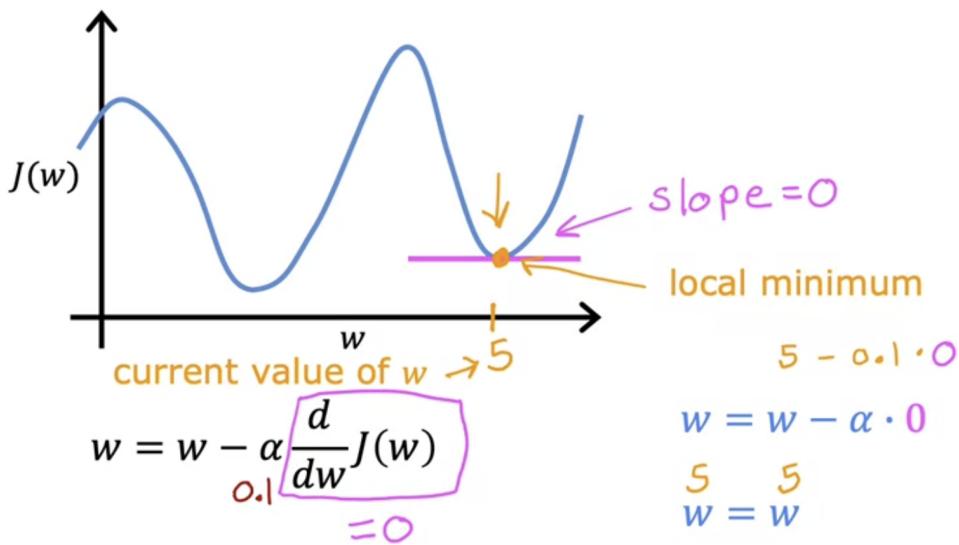
$$w = w - \alpha \frac{d}{dw} J(w)$$

If α is too small...
Gradient descent may be slow.



If α is too large...
Gradient descent may:
- Overshoot, never reach minimum
- Fail to converge, diverge





Linear regression model Cost function

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

Gradient descent algorithm

repeat until convergence {

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

}

Gradient descent algorithm

```

repeat until convergence {
     $w = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$ 
     $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$ 
}

```

Gradient descent algorithm		Assignment	Truth assertion
Repeat until convergence		$a = c$	$a = c$
$\left\{ \begin{array}{l} \underline{w} = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \underline{b} = b - \alpha \frac{\partial}{\partial b} J(w, b) \end{array} \right.$	Learning rate Derivative Simultaneously update w and b	$a = a + 1$ Code	$a = a + 1$ Math a==c
Correct: Simultaneous update $\left. \begin{array}{l} \text{tmp_w} = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \text{tmp_b} = b - \alpha \frac{\partial}{\partial b} J(w, b) \\ w = \text{tmp_w} \\ b = \text{tmp_b} \end{array} \right\}$		Incorrect $\left. \begin{array}{l} \text{tmp_w} = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \text{tmp_b} = b - \alpha \frac{\partial}{\partial b} J(\text{tmp_w}, b) \\ b = \text{tmp_b} \end{array} \right\}$	

- Gradient descent is an algorithm for finding values of parameters w and b that minimize the cost function J.

repeat until convergence {

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

When $\frac{\partial J(w,b)}{\partial w}$ is a negative number (less than zero), what happens to w after one update step?

- w decreases
- w stays the same
- w increases.
- It is not possible to tell if w will increase or decrease.



Correct

The learning rate is always a positive number, so if you take W minus a negative number, you end up with a new value for W that is larger (more positive).

- For linear regression, what is the update step for parameter b?

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



Correct

The update step is $b = b - \alpha \frac{\partial J(w,b)}{\partial b}$ where $\frac{\partial J(w,b)}{\partial b}$ can be computed with this expression:

$$\sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

Gradient descent summary

So far in this course, you have developed a linear model that predicts $f_{w,b}(x^{(i)})$:

$$f_{w,b}(x^{(i)}) = wx^{(i)} + b \quad (1)$$

In linear regression, you utilize input training data to fit the parameters w, b by minimizing a measure of the error between our predictions $f_{w,b}(x^{(i)})$ and the actual data $y^{(i)}$. The measure is called the *cost*, $J(w, b)$. In training you measure the cost over all of our training samples $x^{(i)}, y^{(i)}$

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

In lecture, *gradient descent* was described as:

$$\begin{aligned} &\text{repeat until convergence: } \{ \\ &w = w - \alpha \frac{\partial J(w, b)}{\partial w} \\ &b = b - \alpha \frac{\partial J(w, b)}{\partial b} \\ &\} \end{aligned} \quad (3)$$

where, parameters w, b are updated simultaneously.

The gradient is defined as:

$$\frac{\partial J(w, b)}{\partial w} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{w,b}(x^{(i)}) - y^{(i)})x^{(i)} \quad (4)$$

$$\frac{\partial J(w, b)}{\partial b} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{w,b}(x^{(i)}) - y^{(i)}) \quad (5)$$

Here *simultaneously* means that you calculate the partial derivatives for all the parameters before updating any of the parameters.

Das Gradientenverfahren ist ein allgemeiner Algorithmus zur Optimierung, der optimale Lösungen für eine Vielzahl von Fragestellungen ermitteln kann. Der Grundgedanke beim GV ist, die Parameter iterativ so zu verändern, dass eine Kostenfunktion minimiert wird.

Das Gradientenverfahren berechnet den lokalen Gradienten der Fehlerfunktion in Abhängigkeit vom Parametervektor (ω) und bewegt sich in Richtung eines abfallenden Gradienten. Sobald der Gradient null wird, haben Sie ein Minimum gefunden.

Wenn die Lernrate zu klein ist, muss der Algorithmus viele Iterationen durchlaufen, bevor er konvergiert. Das dauert natürlich sehr lange.

Wenn die Lernrate dagegen zu groß ist, kann es passieren, dass Sie die Täler überspringen und auf der anderen Seite landen, möglicherweise sogar höher als zuvor. Dadurch kann der Algorithmus divergieren, also immer größerer Werte erzeugen und überhaupt keine gute Lösungen finden.

```
#Function to calculate the cost
def compute_cost(x, y, w, b):

    m = x.shape[0]
    cost = 0

    for i in range(m):
        f_wb = w * x[i] + b
        cost = cost + (f_wb - y[i])**2
    total_cost = 1 / (2 * m) * cost

    return total_cost
```

```
def compute_gradient(x, y, w, b):
    """
    Computes the gradient for linear regression
    Args:
        x (ndarray (m,)): Data, m examples
        y (ndarray (m,)): target values
        w,b (scalar)      : model parameters
    Returns
        dj_dw (scalar): The gradient of the cost w.r.t. the parameter
        dj_db (scalar): The gradient of the cost w.r.t. the parameter
    """
    # Number of training examples
    m = x.shape[0]
```

```

dj_dw = 0
dj_db = 0

for i in range(m):
    f_wb = w * x[i] + b
    dj_dw_i = (f_wb - y[i]) * x[i]
    dj_db_i = f_wb - y[i]
    dj_db += dj_db_i
    dj_dw += dj_dw_i
dj_dw = dj_dw / m
dj_db = dj_db / m

return dj_dw, dj_db

```

```
def gradient_descent(x, y, w_in, b_in, alpha, num_iters, cost_fn, gradient_fn):
    """
    """

```

Performs gradient descent to fit w,b. Updates w,b by taking num_iters gradient steps with learning rate alpha

Args:

- x (ndarray (m,)) : Data, m examples
- y (ndarray (m,)) : target values
- w_in,b_in (scalar): initial values of model parameters
- alpha (float): Learning rate
- num_iters (int): number of iterations to run gradient descent
- cost_function: function to call to produce cost
- gradient_function: function to call to produce gradient

Returns:

- w (scalar): Updated value of parameter after running gradient descent
- b (scalar): Updated value of parameter after running gradient descent
- J_history (List): History of cost values
- p_history (list): History of parameters [w,b]

"""

```

# An array to store cost J and w's at each iteration primarily
J_history = []
p_history = []
b = b_in
w = w_in

for i in range(num_iters):
    # Calculate the gradient and update the parameters using
    dj_dw, dj_db = gradient_function(x, y, w , b)

    # Update Parameters using equation (3) above
    b = b - alpha * dj_db
    w = w - alpha * dj_dw

    # Save cost J at each iteration
    if i<1000000:      # prevent resource exhaustion
        J_history.append( cost_function(x, y, w , b))
        p_history.append([w,b])
    # Print cost every at intervals 10 times or as many iterations
    if i% math.ceil(num_iters/10) == 0:
        print(f"Iteration {i:4}: Cost {J_history[-1]:0.2e} "
              f"dj_dw: {dj_dw: 0.3e}, dj_db: {dj_db: 0.3e} "
              f"w: {w: 0.3e}, b:{b: 0.5e}")

return w, b, J_history, p_history #return w and J,w history

```

MSE ist als Kostenfunktion eines linearen Regressionsmodells eine konvexe Funktion. Das bedeutet, wenn Sie zwei beliebige Punkte auf der Kurve auswählen, liegt die lineare Verbindung zwischen diesen beiden niemals unter der Kurve. Das impliziert, dass es keine lokalen Minima gibt, nur ein globales Minimum. Sie ist auch eine stetige Funktion mit einer Steigung, die sich niemals abrupt ändert. Diese zwei Umsände haben eine wichtige Konsequenz: Mit dem Gradientenverfahren kann man sich dem globalen Minimum beliebig annähern (wenn Sie lange genug warten und die Lernrate nicht zu groß ist).

Die Kostenfunktion hat also die Form einer Schüssel. Sind die Merkmale sehr unterschiedlich skaliert, kann es aber eine längliche Schüssel sein.

Sie fragen sich vielleicht, wie man die Anzahl der Epochen bestimmen soll. Wenn diese zu gering ist, werden Sie beim Anhalten des Algorithmus noch immer weit von der optimalen Lösung entfernt sein. Ist sie aber zu hoch, verschwenden Sie Zeit, während sich die Modellparameter nicht mehr verändern. Eine einfache Lösung ist, die Anzahl Iterationen auf einen sehr großen Wert zu setzen, aber den Algorithmus anzuhalten, sobald der Gradientenvektor winzig klein wird — denn das passiert, wenn das Gradientenverfahren das Minimum (beinahe) erreicht hat. Der Gradientenvektor wird als winzig betrachtet, wenn sein Betrag kleiner als eine sehr kleine Zahl ϵ wird, was man auch als Toleranz bezeichnet.

SupervisedMachineLearning-WEEK-1-QUIZ-3

2.4. Multiple Linear Regression

Multiple features (variables)

Size in feet ²	Number of bedrooms	Number of floors	Age of home in years	Price (\$) in \$1000's
x_1	x_2	x_3	x_4	
2104	5	1	45	460
i=2	1416	3	2	232
1534	3	2	30	315
852	2	1	36	178
...

$$j = 1 \dots 4$$

$$n = 4$$

$$x_j = j^{\text{th}} \text{ feature}$$

$$n = \text{number of features}$$

$$\bar{x}^{(i)} = \text{features of } i^{\text{th}} \text{ training example}$$

$$x_j^{(i)} = \text{value of feature } j \text{ in } i^{\text{th}} \text{ training example}$$

$$\bar{x}^{(2)} = [1416 \ 3 \ 2 \ 40]$$

$$x_3^{(2)} = 2$$

Model:

$$\text{Previously: } f_{w,b}(x) = wx + b$$

example

$$f_{w,b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b$$

$$f_{w,b}(\vec{x}) = 0.1 \underset{\substack{\uparrow \\ \text{size}}}{x_1} + 4 \underset{\substack{\uparrow \\ \text{#bedrooms}}}{x_2} + 10 \underset{\substack{\uparrow \\ \text{#floors}}}{x_3} + -2 \underset{\substack{\uparrow \\ \text{years}}}{x_4} + 80 \underset{\substack{\uparrow \\ \text{base price}}}{b}$$

2.4.1. Vectorization

Parameters and features

$$\vec{w} = [w_1 \ w_2 \ w_3] \quad n=3$$

b is a number

$$\vec{x} = [x_1 \ x_2 \ x_3]$$

linear algebra: count from 1

NumPy 

$$w[0] \quad w[1] \quad w[2]$$

```
w = np.array([1.0, 2.5, -3.3])
```

```
b = 4           x[0] x[1] x[2]
```

```
x = np.array([10, 20, 30])
```

code: count from 0

Without vectorization $n=100,000$

$$f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

```
f = w[0] * x[0] +
    w[1] * x[1] +
    w[2] * x[2] + b
```



Without vectorization

$$f_{\vec{w},b}(\vec{x}) = \left(\sum_{j=1}^n w_j x_j \right) + b \quad \sum_{j=1}^n \rightarrow j=1 \dots n$$

range(0,n) $\rightarrow j=0 \dots n-1$

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



Vectorization

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

```
f = np.dot(w,x) + b
```



Codes with vectorization perform fast than codes without vectorization.

2.4.1.1. Vector Dot Product

The dot product multiplies the values in two vectors element-wise and then sums the result. Vector dot product requires the dimensions of the two vectors to be the same.

Gradient descent $\vec{w} = (w_1 \ w_2 \ \dots \ w_{16})$ ~~b~~ parameters
 derivatives $\vec{d} = (d_1 \ d_2 \ \dots \ d_{16})$

```
w = np.array([0.5, 1.3, ... 3.4])
d = np.array([0.3, 0.2, ... 0.4])
```

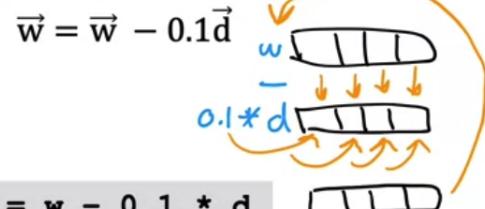
compute $w_j = w_j - 0.1d_j$ for $j = 1 \dots 16$

Without vectorization

$$\begin{aligned} w_1 &= w_1 - 0.1d_1 \\ w_2 &= w_2 - 0.1d_2 \\ &\vdots \\ w_{16} &= w_{16} - 0.1d_{16} \end{aligned}$$

```
for j in range(0,16):
    w[j] = w[j] - 0.1 * d[j]
```

With vectorization



```
w = w - 0.1 * d
```

SupervisedMachineLearning-WEEK-2-QUIZ-1

2.5. Feature Scaling

Feature and parameter values

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

x_1 : size (feet²)
 range: 300 – 2,000 x_2 : # bedrooms
 ↓ ↓
 size #bedrooms large small

House: $x_1 = 2000$, $x_2 = 5$, $\text{price} = \$500k$ one training example

size of the parameters w_1, w_2 ?

$$w_1 = 50, \quad w_2 = 0.1, \quad b = 50$$

$$\begin{aligned}\widehat{\text{price}} &= \frac{50 * 2000}{100,000K} + \frac{0.1 * 5}{0.5K} + \frac{50}{50K} \\ \widehat{\text{price}} &= \$100,050.5K = \$100,050,500\end{aligned}$$

$$w_1 = 0.1, \quad w_2 = 50, \quad b = 50$$

small large

$$\begin{aligned}\widehat{\text{price}} &= \frac{0.1 * 2000K}{200K} + \frac{50 * 5}{250K} + \frac{50}{50K} \\ \widehat{\text{price}} &= \$500K\end{aligned}$$

$\widehat{\text{price}} = \$500K$ more reasonable

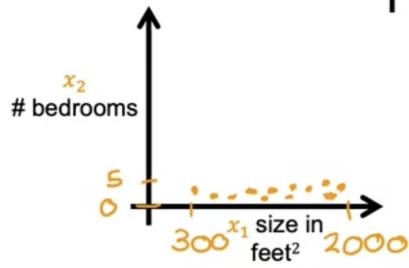
Stanford | ONLINE

DeepLearning.AI

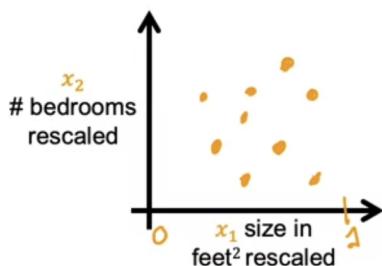
Andrew Ng

- Divide each number by the maximum number.

Feature scaling



$$\begin{aligned}300 \leq x_1 &\leq 2000 & 0 \leq x_2 &\leq 5 \\ x_{1,scaled} &= \frac{x_1}{2000} & x_{2,scaled} &= \frac{x_2}{5} \\ 0.15 \leq x_{1,scaled} &\leq 1 & 0 \leq x_{2,scaled} &\leq 1\end{aligned}$$

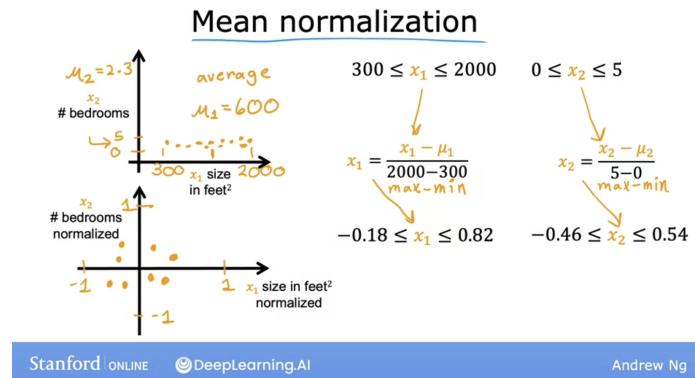


Stanford | ONLINE

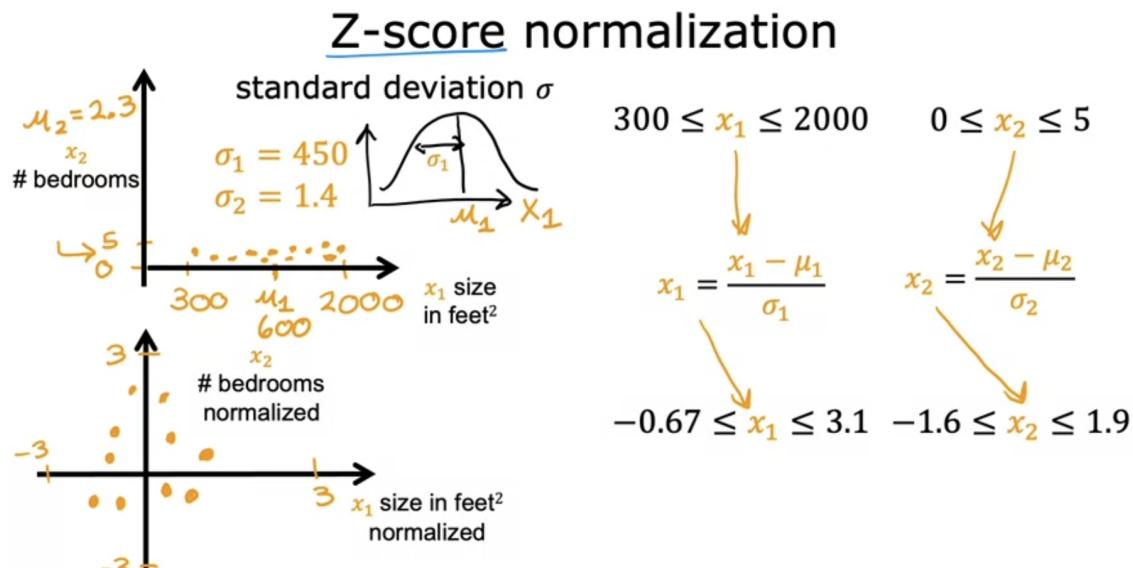
DeepLearning.AI

Andrew Ng

- Mean normalization (Min-max scaling): find the average (mean) and then minus each number by the calculated mean and then divide them by (max - min).



- Z-score normalization: calculate the standard deviation



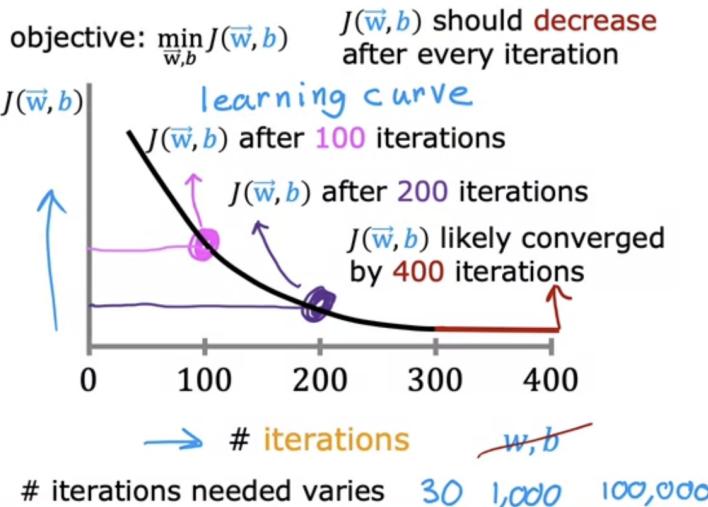
Feature scaling

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

$$\begin{array}{l} -3 \leq x_j \leq 3 \\ -0.3 \leq x_j \leq 0.3 \end{array} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{acceptable ranges}$$

$0 \leq x_1 \leq 3$	Okay, no rescaling
$-2 \leq x_2 \leq 0.5$	Okay, no rescaling
$-100 \leq x_3 \leq 100$	too large → rescale
$-0.001 \leq x_4 \leq 0.001$	too small → rescale
$98.6 \leq x_5 \leq 105$	too large →

Make sure gradient descent is working correctly



Automatic convergence test

Let ϵ "epsilon" be 10^{-3} .
 0.001

If $J(\vec{w}, b)$ decreases by $\leq \epsilon$ in one iteration,
declare convergence.

(found parameters \vec{w}, b
to get close to
global minimum)

2.6. Identify problem with gradient descent

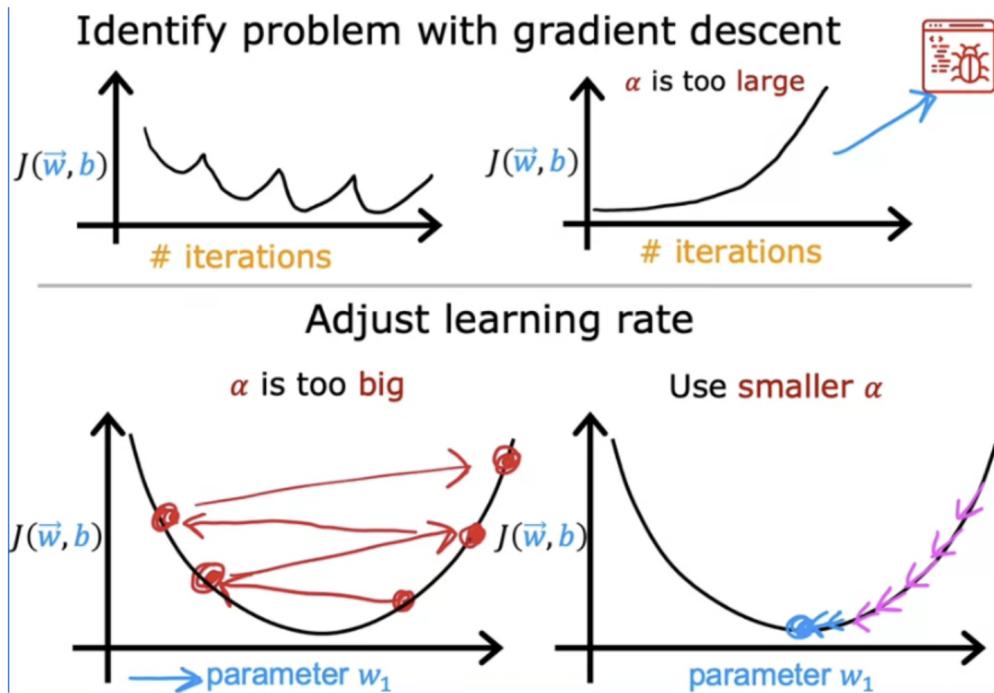
▼ Gradient Descent should consistently decrease the loss function, provided everything is set up correctly. If loss is bouncing up and down, it could be due to one or more of the following reasons:

▼ Learning Rate is too large

A large learning rate can cause the gradient updates to overshoot the optimal point, making the loss oscillate or even diverge.

▼ Bug in the code

Errors in implementing the gradient computation, backpropagation, or updating the weights.



Gradient Descent With Multiple Variables

Here are the equations you developed in the last lab on gradient descent for multiple variables.:

$$\begin{aligned} \text{repeat until convergence: } & \{ \\ w_j &:= w_j - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial w_j} \quad \text{for } j = 0..n-1 \\ b &:= b - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial b} \\ \} \end{aligned} \quad (1)$$

where, n is the number of features, parameters w_j, b , are updated simultaneously and where

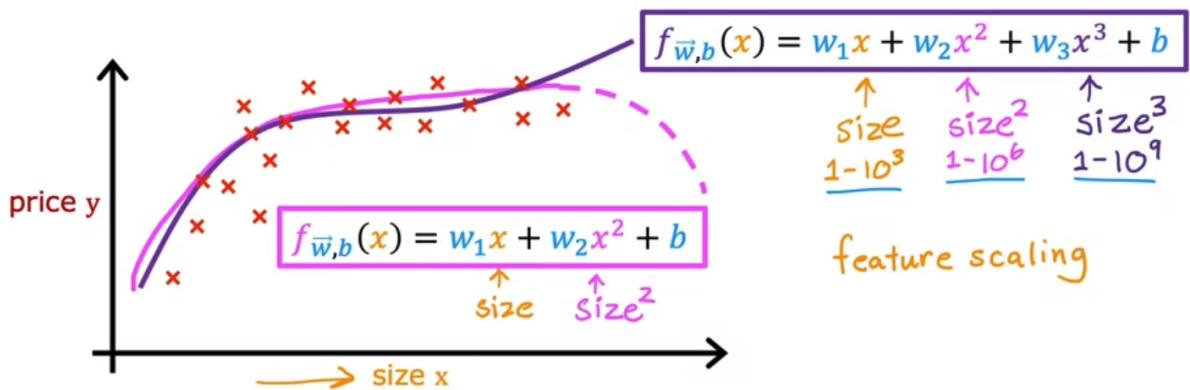
$$\frac{\partial J(\mathbf{w}, b)}{\partial w_j} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)} \quad (2)$$

$$\frac{\partial J(\mathbf{w}, b)}{\partial b} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - y^{(i)}) \quad (3)$$

- m is the number of training examples in the data set
- $f_{\mathbf{w}, b}(\mathbf{x}^{(i)})$ is the model's prediction, while $y^{(i)}$ is the target value

2.7. Polynomial Regression

Polynomial regression



SupervisedMachineLearning-WEEK-2-QUIZ-2

III. Feature engineering

Feature engineering: Using intuition to design new features, by transforming or combining original features.

Feature engineering

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

frontage depth

$$\text{area} = \text{frontage} \times \text{depth}$$

$$x_3 = x_1 x_2$$

new feature

$$f_{\vec{w}, b}(\vec{x}) = \underline{w_1} \underline{x_1} + \underline{w_2} \underline{x_2} + \underline{w_3} x_3 + b$$



Feature engineering:
Using **intuition** to design
new features, by
transforming or combining
original features.

Stanford | ONLINE

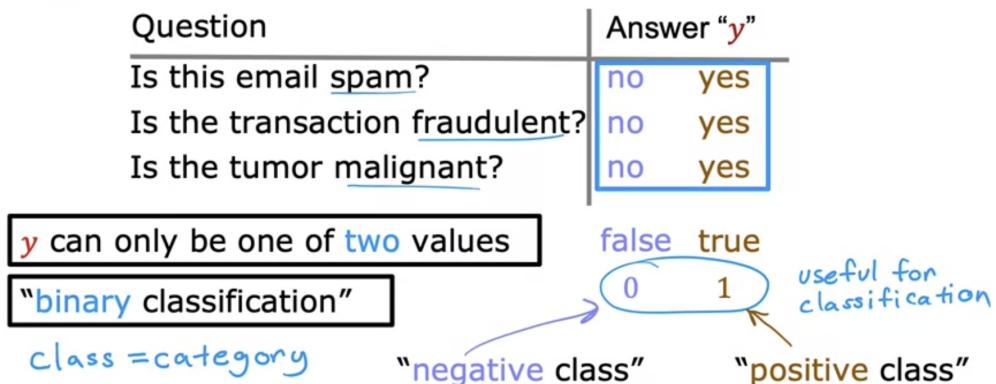
DeepLearning.AI

Andrew Ng

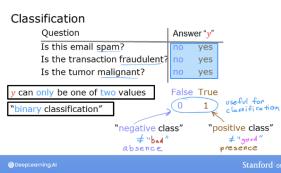
WEEK2_ASSIGNMENT_VADHA_SAMEDY_HUN.ipynb

4. Classification (Logistic Regression)

Classification



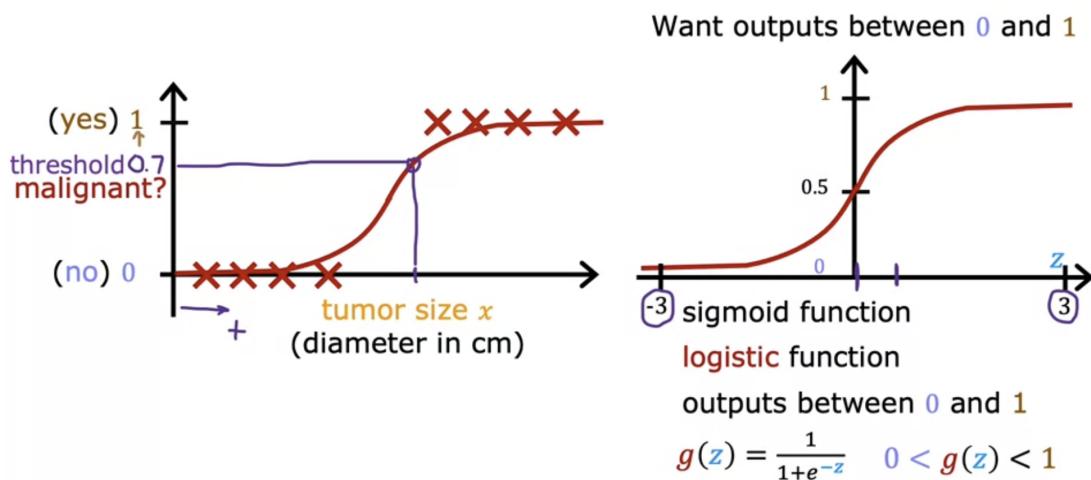
Classification Problems



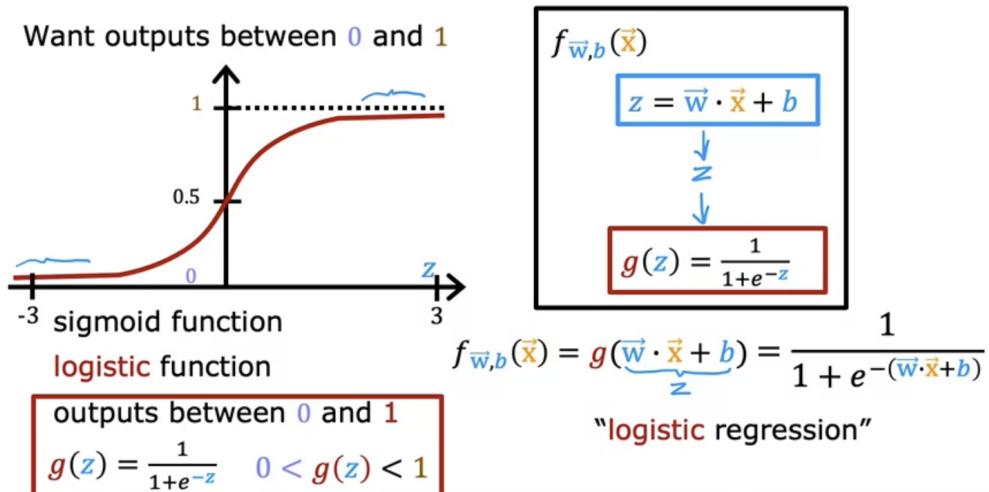
Examples of classification problems are things like: identifying email as Spam or Not Spam or determining if a tumor is malignant or benign. In particular, these are examples of *binary classification* where there are two possible outcomes. Outcomes can be described in pairs of 'positive'/negative' such as 'yes'/no, 'true'/false' or '1'/0'.

Plots of classification data sets often use symbols to indicate the outcome of an example. In the plots below, 'X' is used to represent the positive values while 'O' represents negative outcomes.

4.1. Logistic regression

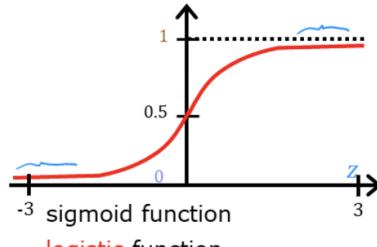


Sigmoid function (logistic function) → output 0 and 1.



Sigmoid or Logistic Function

Want outputs between 0 and 1



outputs between 0 and 1

$$g(z) = \frac{1}{1+e^{-z}} \quad 0 < g(z) < 1$$

DeepLearning.AI

As discussed in the lecture videos, for a classification task, we can start by using our linear regression model, $f_{w,b}(x^{(i)}) = w \cdot x^{(i)} + b$, to predict y given x .

- However, we would like the predictions of our classification model to be between 0 and 1 since our output variable y is either 0 or 1.
- This can be accomplished by using a "sigmoid function" which maps all input values to values between 0 and 1.

Let's implement the sigmoid function and see this for ourselves.

Formula for Sigmoid function

The formula for a sigmoid function is as follows -

$$g(z) = \frac{1}{1+e^{-z}} \quad (1)$$

In the case of logistic regression, z (the input to the sigmoid function), is the output of a linear regression model.

- In the case of a single example, z is scalar.
- in the case of multiple examples, z may be a vector consisting of m values, one for each example.
- The implementation of the sigmoid function should cover both of these potential input formats. Let's implement this in Python.

Logistic Regression: a logistic regression model applies the sigmoid to the familiar linear regression model as shown:

Logistic Regression

A logistic regression model applies the sigmoid to the familiar linear regression model as shown below:

$$f_{w,b}(x^{(i)}) = g(w \cdot x^{(i)} + b) \quad (2)$$

where

$$g(z) = \frac{1}{1+e^{-z}} \quad (3)$$

$$f_{w,b}(x^{(i)}) = g(\underbrace{w \cdot x^{(i)} + b}_z) = \frac{1}{1 + e^{-(w \cdot x^{(i)} + b)}}$$

"logistic regression"

Stanford ONLINE

Squared error cost

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

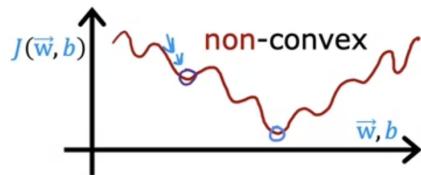
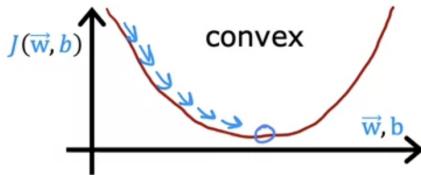
loss $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$

linear regression

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

logistic regression

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$



Stanford ONLINE

DeepLearning.AI

Andrew Ng

Squared error for logistic regression?

Squared Error Cost

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

average of training set

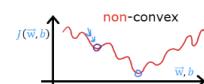
$$L(y^{(i)}, y^{(i)}) = \frac{1}{2} (y^{(i)} - f_{\vec{w}, b}(\vec{x}^{(i)}))^2$$

single training example

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, f_{\vec{w}, b}(\vec{x}^{(i)}))$$

linear regression $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

$$\text{logistic regression } f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$



Recall for Linear Regression we have used the **squared error cost function**: The equation for the squared error cost with one variable is:

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

where

$$f_{w, b}(x^{(i)}) = w x^{(i)} + b \quad (2)$$

DeepLearning.AI

Stanford ONLINE

For the simplified loss function:

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

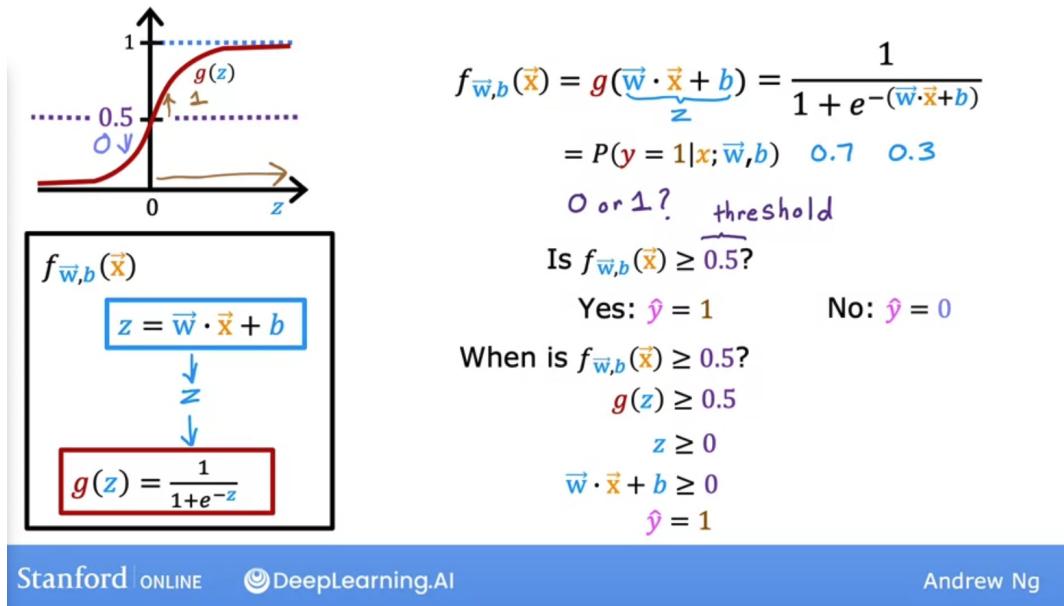
if the target $y^{(i)} = 1$, then what does this expression simplify to?

- $-\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$
- $-\log(f_{\vec{w}, b}(\vec{x}^{(i)}))$

✓ Correct

The second term of the expression is reduced to zero when the target equals 1.

4.1.1. Decision boundary



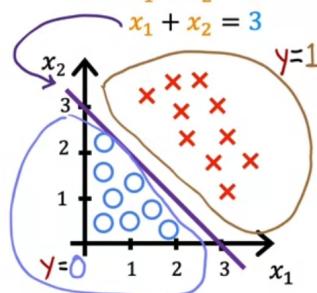
Decision boundary

$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 x_1 + w_2 x_2 + b)$$

$$\text{Decision boundary } z = \vec{w} \cdot \vec{x} + b = 0$$

$$z = x_1 + x_2 - 3 = 0$$

$$x_1 + x_2 = 3$$



4.1.2. Cost function

Training set

	tumor size (cm) x_1	...	patient's age x_n	malignant? y	$i = 1, \dots, m \leftarrow$ training examples $j = 1, \dots, n \leftarrow$ features
$i=1$	10		52	1	
:	2		73	0	
:	5		55	0	
$i=m$	12		49	1	$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$
	

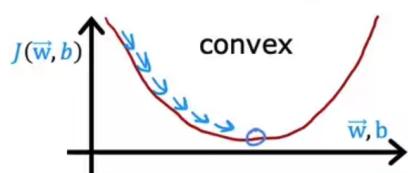
Squared error cost

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

loss $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$

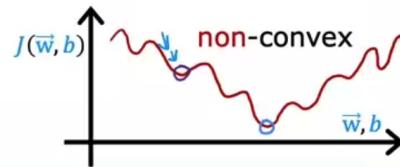
linear regression

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$



logistic regression

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$



Cost

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$$

$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$ loss
 $\begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$
Convex
can reach a global minimum

SupervisedMachineLearning-WEEK-3-QUIZ-2

4.2. Gradient descent for logistic regression

Gradient descent for logistic regression

```

repeat {           looks like linear regression!
     $w_j = w_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right]$ 
     $b = b - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \right]$ 
} simultaneous updates

```

- Same concepts:
- Monitor gradient descent (learning curve)

Linear regression $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

Logistic regression $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

Gradient descent for logistic regression

repeat {

looks like linear regression!

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

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

} simultaneous updates

Same concepts:

- Monitor gradient descent (learning curve)
- Vectorized implementation
- Feature scaling

Linear regression $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

Logistic regression $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

Stanford ONLINE

DeepLearning.AI

Andrew Ng

Logistic Gradient Descent

Recall the gradient descent algorithm utilizes the gradient calculation:

repeat until convergence: {

$$\begin{aligned} w_j &= w_j - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial w_j} && \text{for } j := 0..n-1 \\ b &= b - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial b} \\ \} \end{aligned}$$

Where each iteration performs simultaneous updates on w_j for all j , where

$$\frac{\partial J(\mathbf{w}, b)}{\partial w_j} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

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

- m is the number of training examples in the data set

- $f_{\mathbf{w}, b}(\mathbf{x}^{(i)})$ is the model's prediction, while $y^{(i)}$ is the target

- For a logistic regression model

$$z = \mathbf{w} \cdot \mathbf{x} + b$$

$$f_{\mathbf{w}, b}(x) = g(z)$$

where $g(z)$ is the sigmoid function:

$$g(z) = \frac{1}{1+e^{-z}}$$

Gradient descent for logistic regression

repeat {

looks like linear regression!

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

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

} simultaneous updates

Linear regression $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

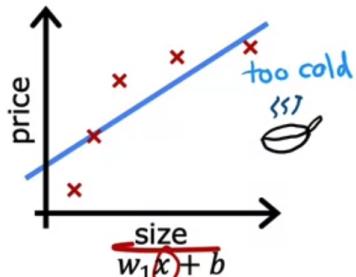
Logistic regression $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

Andrew Ng

SupervisedMachineLearning-WEEK-3-QUIZ-3

4.3. Overfitting

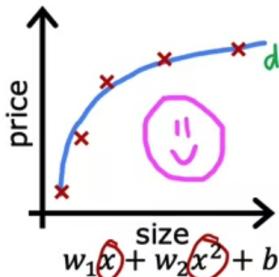
Regression example



underfit

- Does not fit the training set well

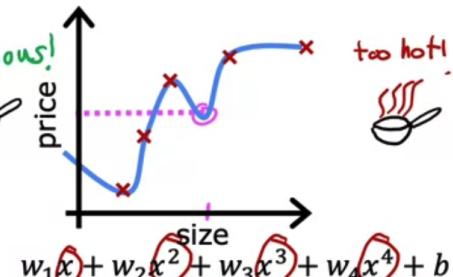
high bias



just right

- Fits training set pretty well

generalization

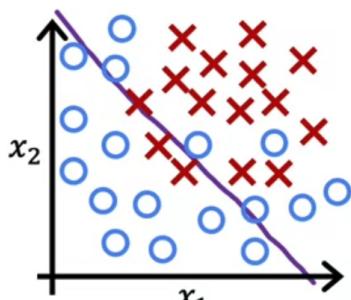


overfit

- Fits the training set extremely well

high variance

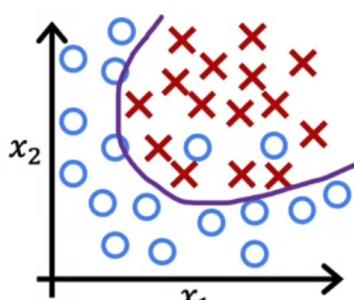
Classification



$$z = w_1x_1 + w_2x_2 + b$$

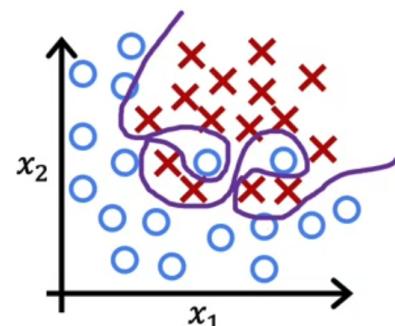
$$f_{\vec{w}, b}(\vec{x}) = g(z)$$

g is the sigmoid function
underfit high bias



$$\begin{aligned} z &= w_1x_1 + w_2x_2 \\ &+ w_3x_1^2 + w_4x_2^2 \\ &+ w_5x_1x_2 + b \end{aligned}$$

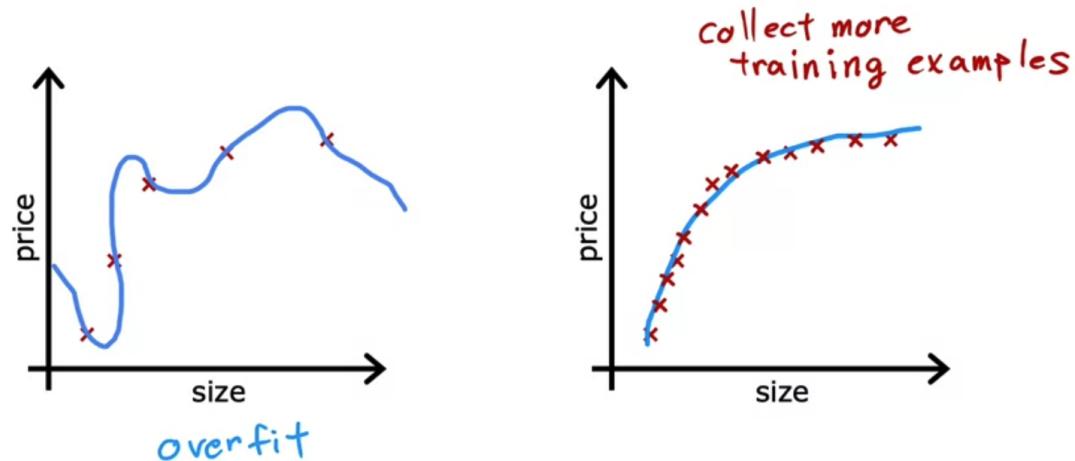
just right



$$\begin{aligned} z &= w_1x_1 + w_2x_2 \\ &+ w_3x_1^2x_2 + w_4x_1^2x_2^2 \\ &+ w_5x_1^2x_2^3 + w_6x_1^3x_2 \\ &+ \dots + b \end{aligned}$$

4.3.1. More data to reduce Overfitting

Collect more training examples



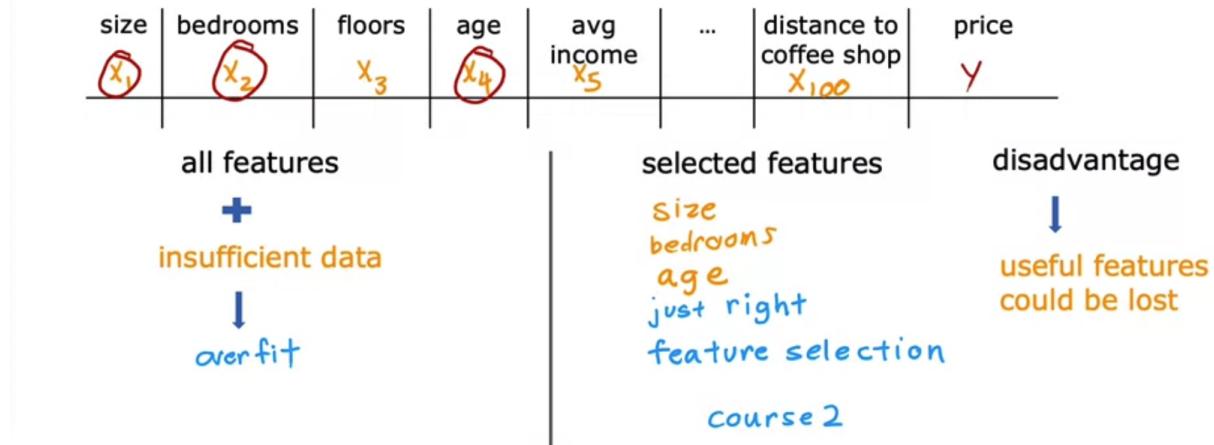
Stanford | ONLINE

DeepLearning.AI

Andrew Ng

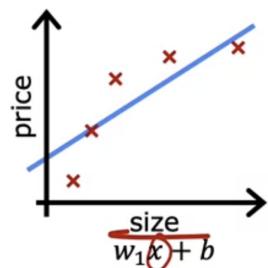
4.3.2. Select feature(s) to include/exclude

Select features to include/exclude



4.3.3. Regularization to Reduce Overfitting

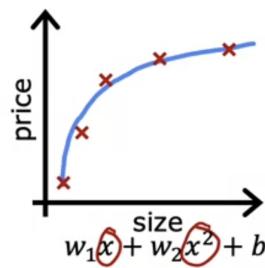
Regression example



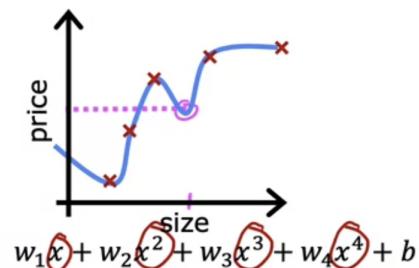
underfit

- Does not fit the training set well

high bias



generalization



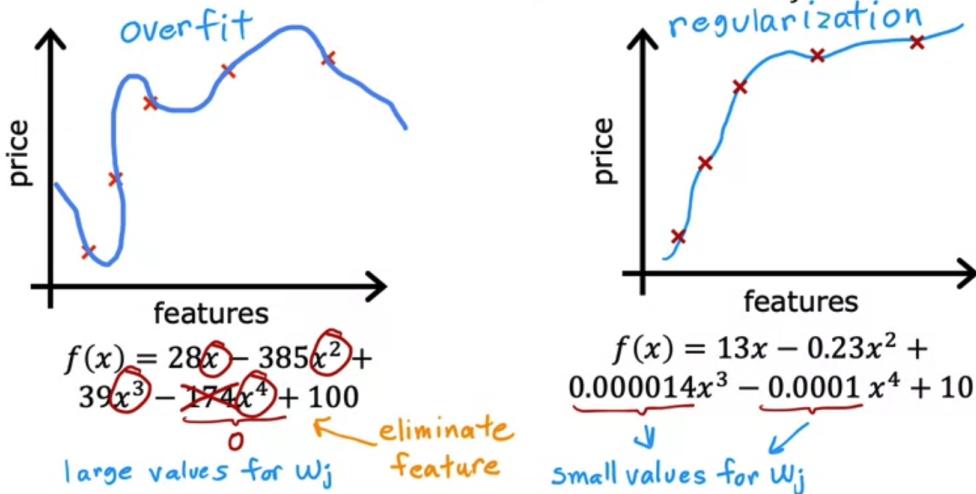
overfit

- Fits the training set extremely well

high variance

Regularization

Reduce the size of parameters w_j

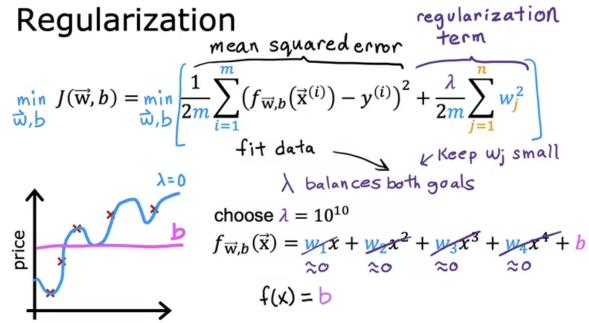


4.3.4. Solutions

To reduce overfitting:

- Gather more data.
- Select features to include/exclude (just right, feature selection)
- Regularization (penalized): reduce the size of parameters w_j (either eliminate a feature (that has big impact) or give small values)

4.3.5. Regularization



How to choose lambda?

If lambda too large \rightarrow underfit.

If lambda too small \rightarrow overfit.

Lambda value needs to appropriately balance the trade-off of minimizing the MSE and keeping the parameter small.

For a model that includes the regularization parameter λ (lambda), increasing λ will tend to...

- Increase the size of parameter b .
- Decrease the size of parameters w_1, w_2, \dots, w_n .
- Increases the size of the parameters w_1, w_2, \dots, w_n
- Decrease the size of the parameter b .

Correct

Increasing the regularization parameter *lambda* reduces overfitting by reducing the size of the parameters. For some parameters that are near zero, this reduces the effect of the associated features.

Implementing gradient descent

```

repeat {
     $w_j = w_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m [(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}] + \frac{\lambda}{m} w_j \right]$ 
     $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$ 
} simultaneous update  $j = 1, \dots, n$ 
 $w_j = \underbrace{w_j - \alpha \frac{\lambda}{m} w_j}_{w_j \left(1 - \alpha \frac{\lambda}{m}\right)} - \underbrace{\alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}}_{\text{usual update}}$ 

```

$\alpha \frac{\lambda}{m}$
 $0.01 \frac{1}{50} = 0.0002$
 $w_j (1 - 0.0002)$
 0.9998

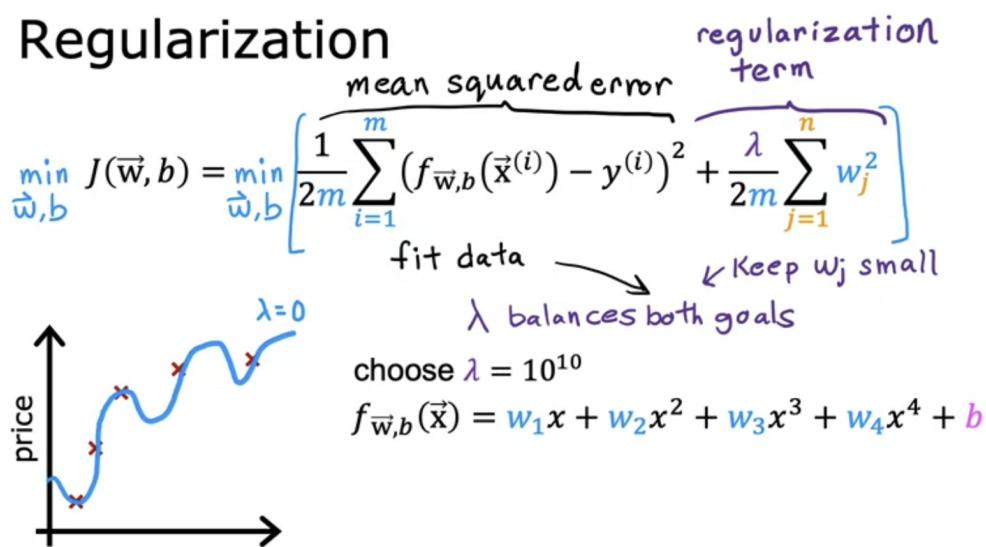
For regularized **logistic** regression, how do the gradient descent update steps compare to the steps for linear regression?

- They look very similar, but the $f(x)$ is not the same.
- They are identical

Correct

For logistic regression, $f(x)$ is the sigmoid (logistic) function, whereas for linear regression, $f(x)$ is a linear function.

Regularization



Implementing gradient descent

repeat {

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

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

} simultaneous update $j = 1 \dots n$

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

$\alpha \frac{\lambda}{m}$
 $0.01 \frac{1}{50} = 0.0002$
 $w_j \left(1 - 0.0002 \right)$
 0.9998

How we get the derivative term (optional)

$$\begin{aligned}
 \frac{\partial}{\partial w_j} J(\vec{w}, b) &= \frac{\partial}{\partial w_j} \left[\frac{1}{2m} \sum_{i=1}^m \underbrace{(f(\vec{x}^{(i)}) - y^{(i)})^2}_{\vec{w} \cdot \vec{x}^{(i)} + b} + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right] \\
 &= \frac{1}{2m} \sum_{i=1}^m \left[(\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)}) \cancel{\cancel{x}_j^{(i)}} \right] + \frac{\lambda}{2m} \cancel{\cancel{2}} w_j \text{ No } \sum_{j=1}^n \\
 &= \frac{1}{m} \sum_{i=1}^m \left[(\underbrace{\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)}}_{f(\vec{x})}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j \\
 &= \frac{1}{m} \sum_{i=1}^m \left[(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j
 \end{aligned}$$

- Regularization penalizes “memorization” (over-learning examples).
- Helps the model generalize to unseen examples.
- Changes the representations of learning (either more sparse or more distributed depending on the regularizer).
- Can increase or decrease training time.
- Can decrease training accuracy but increase generalization.
- Works better for large models with multiple hidden layers.
- Generally works better with sufficient data.

▼ How to think about regularization?

- Adds a cost to the complexity of the solution.
- Forces the solution to be smooth.
- Prevents the model from learning item-specific details.

4.4. Types

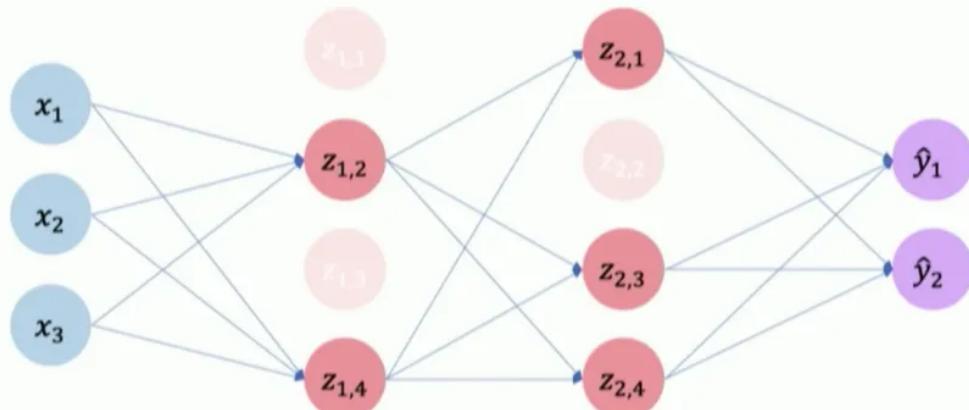
- Modify the model (dropout) → "Node regularization"
- Add a cost to the loss function (L1/2) → "Loss regularization"
- Modify or add data (batch training, data augmentation) → "Data regularization"

4.4.1.1. Dropout

"Remove" nodes randomly during learning (force activation = 0).

Regularization I: Dropout

- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node



▼ EFFECT?

- Prevents a single node from learning too much.
- Forces the model to have distributed representations.
- Makes the model less reliant on individual nodes and thus more stable.

- Generally requires more training (though each epoch computes faster).
- Can decrease training accuracy but increase generalization.
- Usually works better on deep than shallow networks.
- Debate about applying it to convolution layers.
- Works better with sufficient data, unnecessary with “enough” data.

4.4.1.2. L1/L2 regularization

“weight decay” add a cost to the loss function to prevent weights from getting too large.

Loss function
without regularization

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})$$

Loss function
with regularization

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) + \lambda \sum_{j=1}^K \|\theta_j^2\|$$

Penalty term

- L1 regularization (“lasso”) creates sparse weights by setting some to 0. (sparse weights (lots of 0's))

- L2 regularization ("ridge" or "weight decay") shrinks all weights, especially large weights. (>0 in same range)

▼ WHEN TO REGULARIZE?

Regularize as little as possible and as much as necessary.

▼ Other possibilities:

- L1 + L2 ("elastic net" regression)
- Norm of weight matrix
- Sample-specific (positive bias on cancer diagnosis)

▼ Why does regularization reduce overfitting?

- Discourages complex and sample-specific representations.
- Prevents overfitting to training examples.
- Large weights lead to instability (very different outputs for similar inputs).
- In large, complex models with lots of weights (high risk of overfitting).
- Use L1 when trying to understand the important encoding features (more common in regression than DL).
- When training accuracy is much higher than validation accuracy.

4.4.1.3. Data augmentation

"Data augmentation" add more data as slightly modified versions of existing data (usually just for images).

▼ Which regularization method to use?

Oftentimes the "best" method is problem- or architecture-specific.

In many cases, different regularization methods work equally well.

4.5. Training vs. evaluation mode

Gradients are computed only during backprop, not during evaluation.

Some regularization methods are applied only during training, not during evaluation.

Ergo: We need a way to deactivate gradient computations and regularization while evaluating model performance.

4.6. Training in mini-batches

- Batch size is often powers-of-2 (e.g., $2^4 = 16$), between 2 and 512.
- Training in batches can decrease computation time because of vectorization (matrix multiplication instead of for-loops).
- But batching can increase computation time for large batches and large data samples (e.g., images).
- Batching is a form of regularization: It smooths learning by averaging the loss over many samples, and thereby reduces overfitting.

<https://github.com/SamedyHUNX/supervised-learning>