



Classification

Motivations

Classification

Question

Answer “*y*”

Is this email spam?

no yes

Is the transaction fraudulent?

no yes

Is the tumor malignant?

no yes

y can only be one of two values

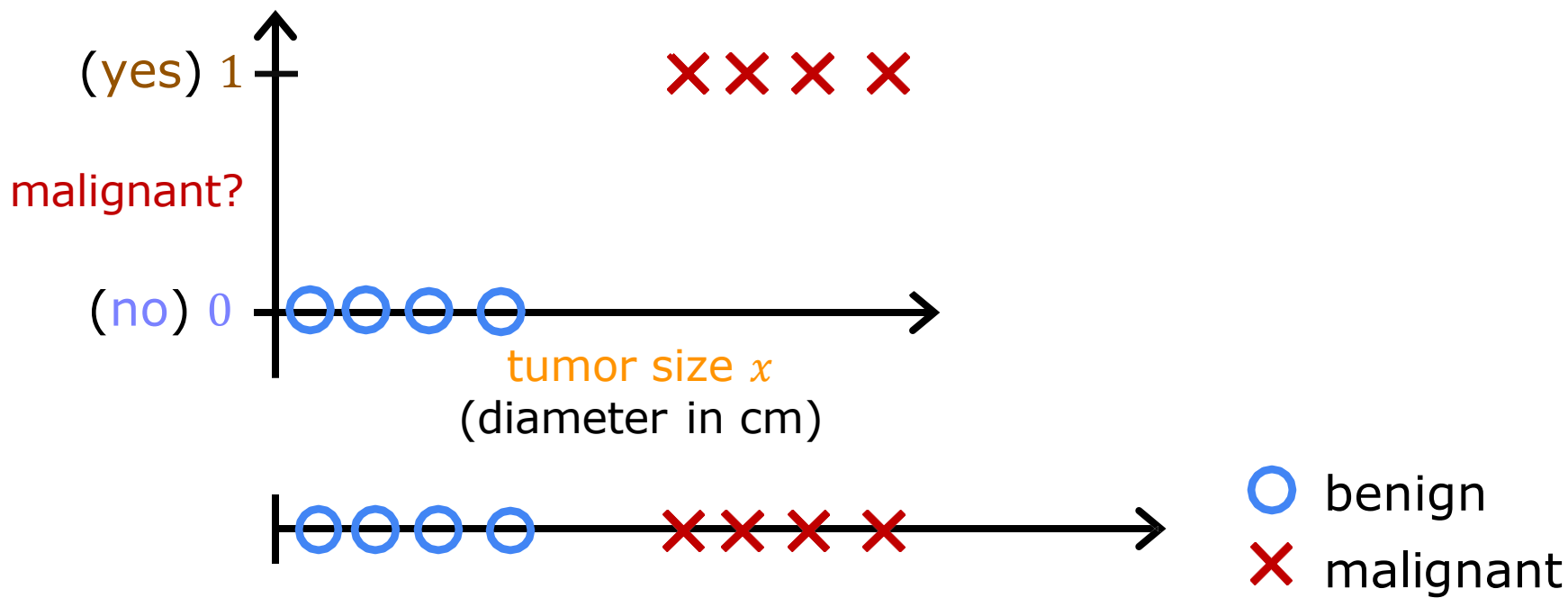
false true

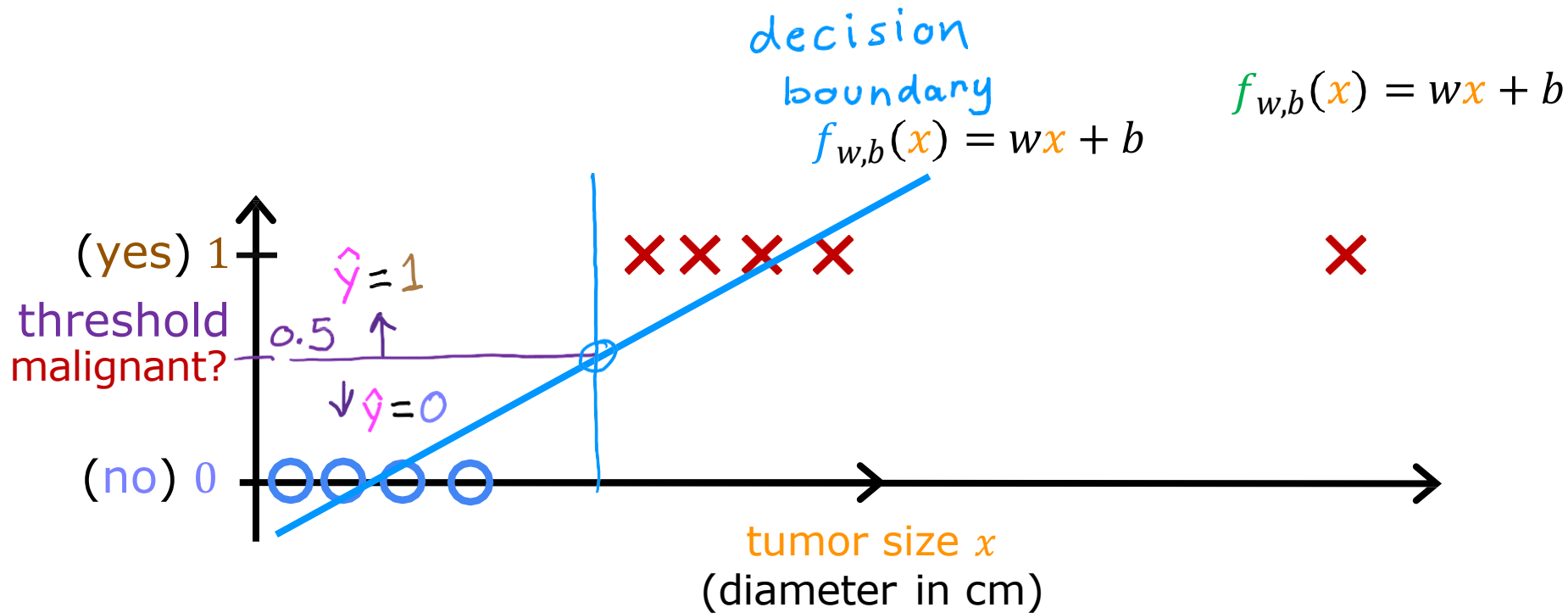
“binary classification”

0

1

class = category





if $f_{w,b}()$

$x < 0.5 \rightarrow y = 0$

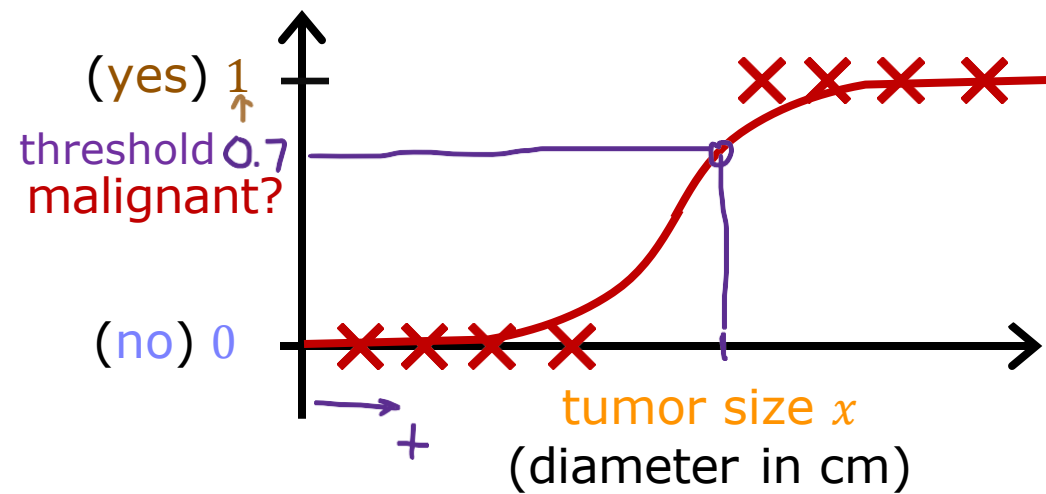
if $f_{w,b}()$

$x \geq 0.5 \rightarrow y = 1$

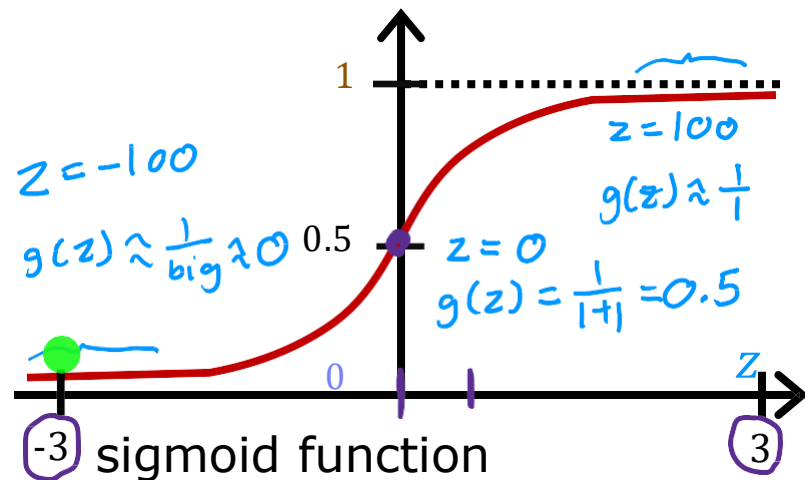


Classification

Logistic Regression



Want outputs between 0 and 1

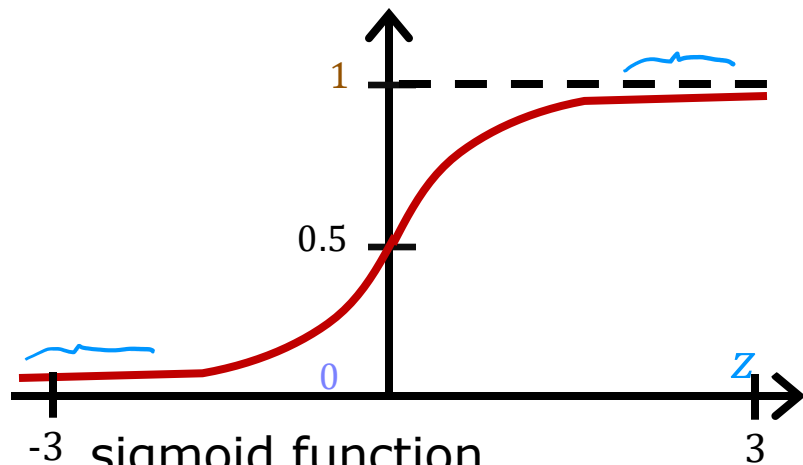


logistic function

outputs between 0 and 1

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

Want outputs between 0 and 1



sigmoid function

logistic function

outputs between 0 and 1

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

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

$$\vec{w} \cdot \vec{x} + b$$



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

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

"logistic regression"

$e \approx 2.7$

Interpretation of logistic regression output

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

“probability” that class is 1

Example:

x is “tumor size”

y is 0 (not malignant)

or 1 (malignant)

$$f_{\vec{w},b}(\vec{x}) = 0.7$$

70% chance that y is 1

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

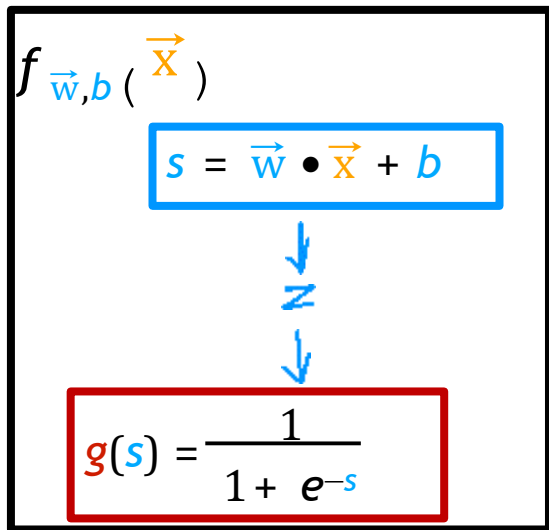
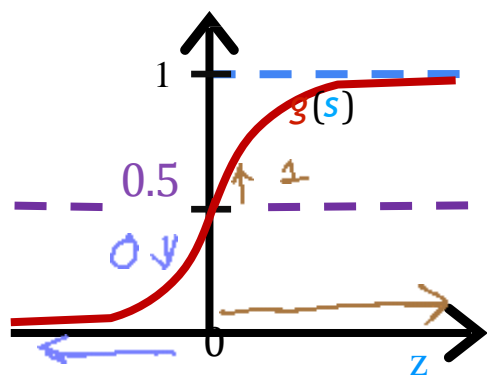
Probability that y is 1,
given input \vec{x} , parameters \vec{w}, b

$$P(y = 0) + P(y = 1) = 1$$



Classification

Decision Boundary



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

$$= P(y = 1 | x; \vec{w}, b) \quad 0.7 \quad 0.3$$

0 or 1? threshold

Is $f_{\vec{w},b}(\vec{x}) \geq 0.5$?

Yes: $\hat{y} = 1$

No: $\hat{y} = 0$

When is

$$f_{\vec{w},b}(\vec{x}) \geq 0.5 \iff g(s) \geq 0.5$$

$$s \geq 0$$

$$\vec{w} \cdot \vec{x} + b \geq 0$$

$$\hat{y} = 1$$

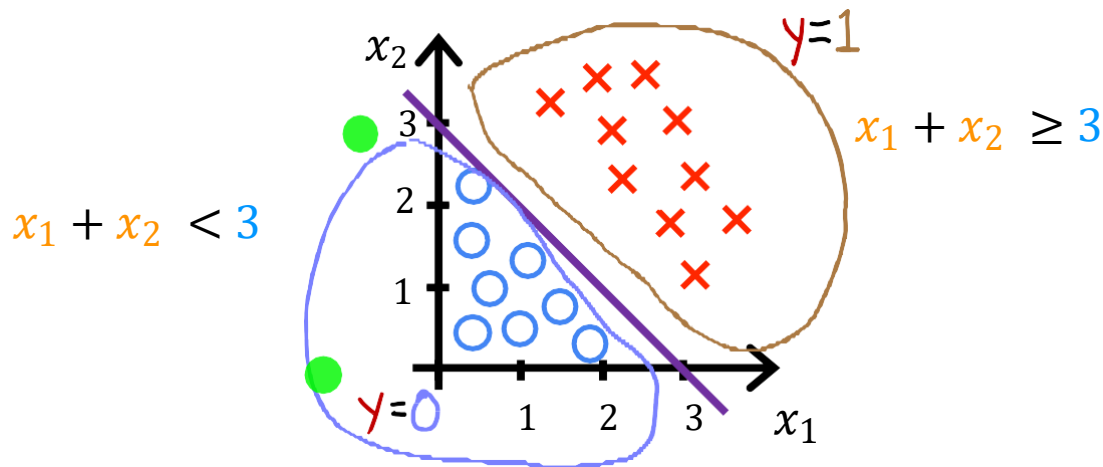
$$\vec{w} \cdot \vec{x} + b < 0$$

$$\hat{y} = 0$$

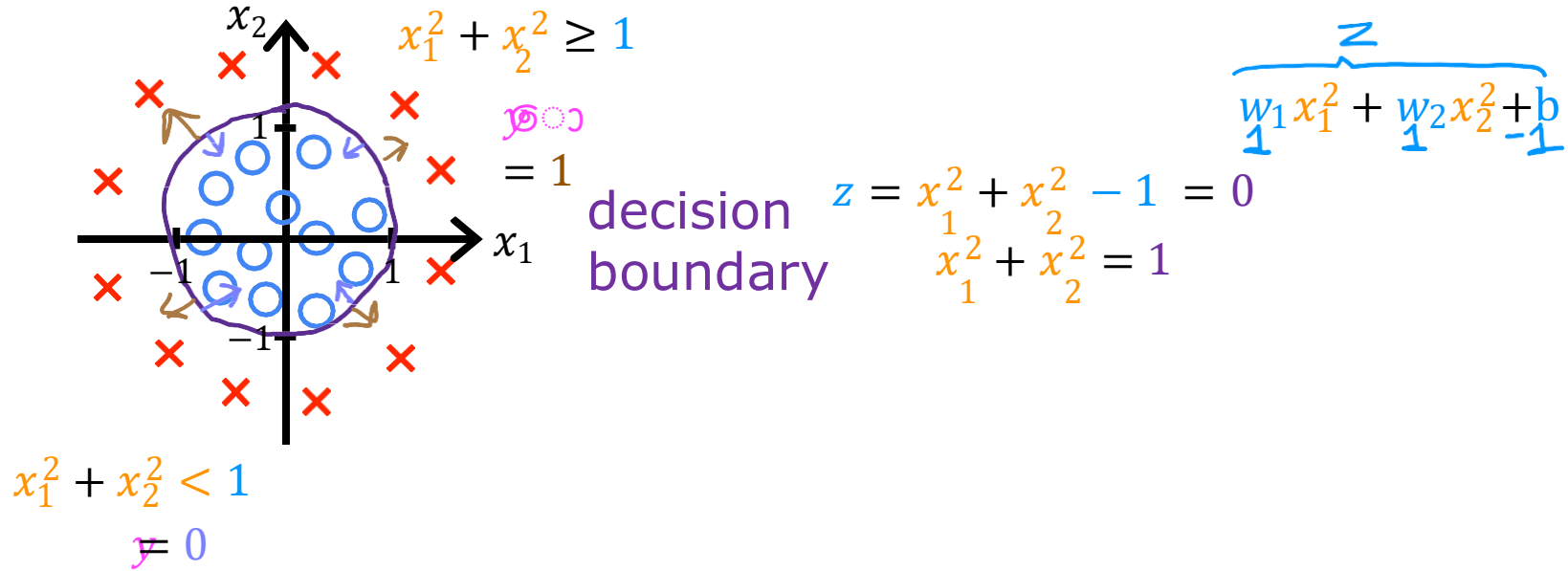
Decision boundary

$$f_{\vec{w},b}(\vec{x}) = g(z) = g(\underbrace{w_1x_1 + w_2x_2 + b}_{\substack{1 \quad 1 \quad -3}})$$

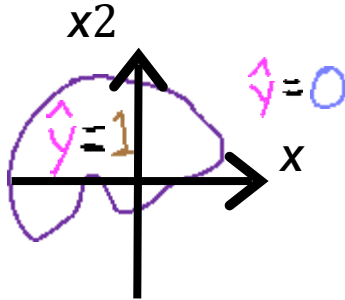
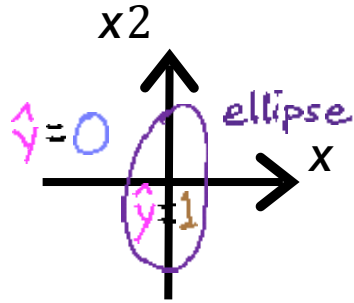
Decision boundary $z = \vec{w} \cdot \vec{x} + b = 0$
 $z = x_1 + x_2 - 3 = 0$
 $x_1 + x_2 = 3$



Non-linear decision boundaries



Non-linear decision boundaries



$$f_{\vec{w}, b}(\vec{x}) = g(s) = g(r_1 x_1 + r_2 x_2 + r_3 x_1^2 + r_4 x_1 x_2 + r_5 x_2^2 + r_6 x_1^3 + \dots + b)$$



Cost Function

Cost Function for
Logistic Regression

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	<div style="border: 1px solid red; padding: 5px; display: inline-block;">target y is 0 or 1</div> $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$
\vdots	2		73	0	
\vdots	5		55	0	
\vdots	12		49	1	
$i=m$	

How to choose $\vec{w} = [w_1 \ w_2 \ \cdots \ w_n]$ and b ?

Squared error cost

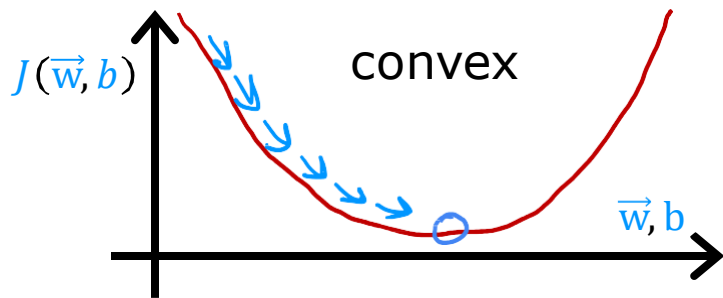
$$\text{cost}$$
$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$$

average of training set

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

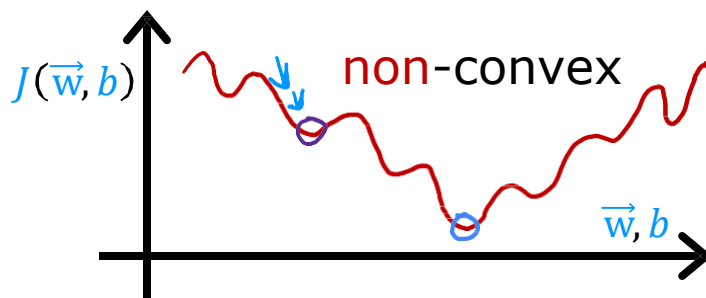
linear regression

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



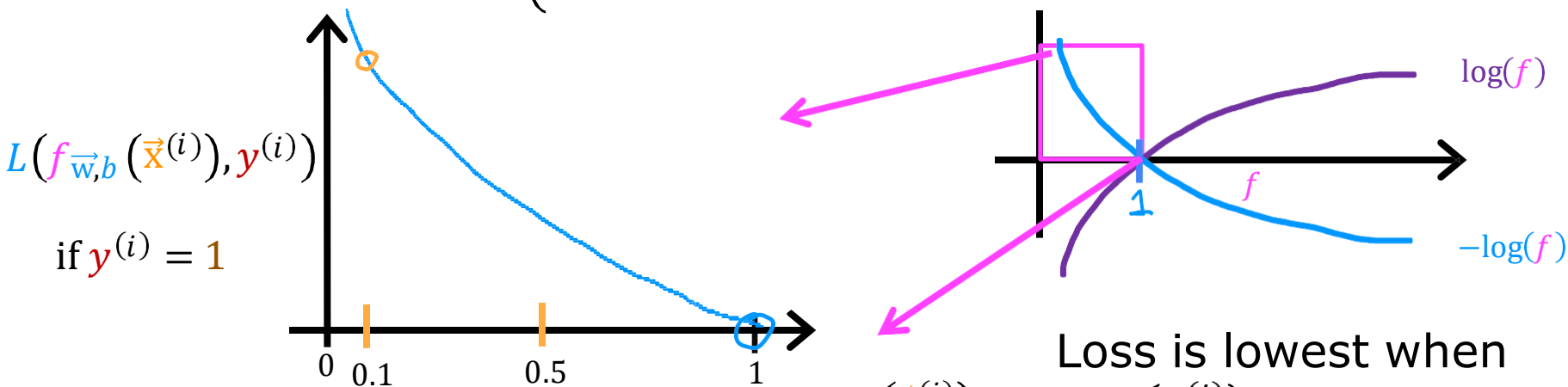
logistic regression

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



Logistic loss function

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



As $f_{\vec{w},b}(\vec{x}^{(i)}) \rightarrow 1$ then loss $\rightarrow 0$

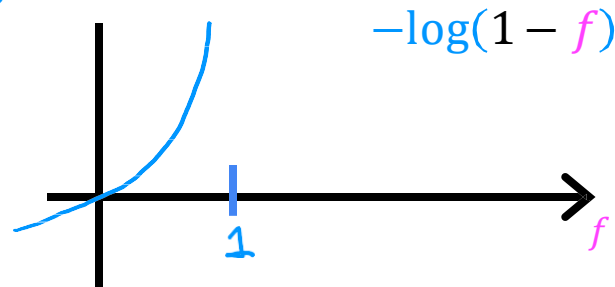
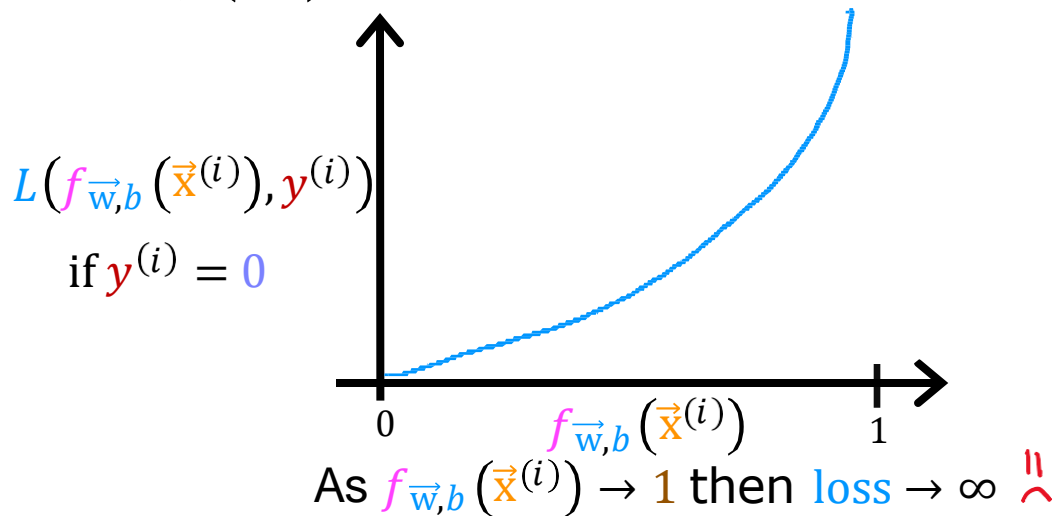
As $f_{\vec{w},b}(\vec{x}^{(i)}) \rightarrow 0$ then loss $\rightarrow \infty$

Loss is lowest when $f_{\vec{w},b}(\vec{x}^{(i)})$ predicts close to true label $y^{(i)}$.

Logistic loss function

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

As $f_{\vec{w},b}(\vec{x}^{(i)}) \rightarrow 0$ then loss $\rightarrow 0$ \Downarrow



The further prediction $f_{\vec{w},b}(\vec{x}^{(i)})$ is from target $y^{(i)}$, the higher the loss.

Cost

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

$$= \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

convex
→ can reach a global minimum

find w, b that minimize cost J



Cost Function

Simplified Cost
Function for Logistic
Regression

Simplified loss function

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

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

if $y^{(i)} = 1$:

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

$\underbrace{(1-1)}_0$

Simplified loss function

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

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

if $y^{(i)} = 1$:

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

if $y^{(i)} = 0$:

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

Simplified cost function

loss

$$L(f_{\vec{w},b}(\vec{x}^{(i)}), y^{(i)}) = \underbrace{-y^{(i)}\log(f_{\vec{w},b}(\vec{x}^{(i)})) - (1 - y^{(i)})\log(1 - f_{\vec{w},b}(\vec{x}^{(i)}))}_{\text{convex (single global minimum)}}$$

cost

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

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

maximum likelihood



Gradient Descent

Gradient Descent Implementation

Training logistic regression

Find \vec{w}, b

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

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

Gradient descent

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

repeat {

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

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

} simultaneous updates

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

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

Gradient descent for logistic regression

repeat {

looks like linear regression!

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

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

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

Logistic regression hypothesis:

$$h_{\theta}(x) = g(\theta^T x)$$

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

Cost function in logistic regression is:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(h_{\theta}(x^i)) + (1 - y^i) \log(1 - h_{\theta}(x^i))]$$

Vectorized implementation:

$$h = g(X\theta)$$

$$J(\theta) = \frac{1}{m} (-y^T \log(h) - (1 - y)^T \log(1 - h))$$

The gradient of the cost is a vector of the same length as θ where j^{th} element (for $j = 0, 1, \dots, n$) is defined as follows:

$$\nabla J(\theta) = \frac{1}{m} \sum_{i=1}^m ((h_{\theta}(x^i) - y^i) \cdot x_j^i)$$

Newton's Method for parameter optimization:

Advantages :

Fast Convergence: Newton's method converges faster than gradient descent, especially for convex and well-behaved functions.

It typically achieves **quadratic convergence** near the optimum, meaning the error shrinks exponentially with each iteration.

Adaptability to Curvature: By using the Hessian, Newton's method adjusts the step size based on the curvature of the function, preventing overshooting or undershooting the minimum.

Precise Steps: The inclusion of second-order information allows for more precise updates, especially in regions where the function's curvature changes rapidly.

Newton's Method for Parameter Optimization:

- **Hessian Calculation:** Computing the Hessian matrix can be expensive, especially for functions with many parameters (large-scale optimization problems).
The Hessian is an $n \times n$ matrix, where n is the number of parameters.
- **Hessian Inversion:** Newton's method requires inverting the Hessian matrix, which can be computationally expensive and numerically unstable if the Hessian is not well-conditioned (i.e., close to singular, i.e. eigenvalues are very small, or the determinant is close to zero).
- **Local Minimum:** Like gradient descent, Newton's method can get stuck in local minima if the function is not convex

Modified or Truncated Newton Methods

- Due to the computational cost of calculating and inverting the Hessian matrix, several modified versions of Newton's method are commonly used in practice.
- These methods make Newton's method more efficient for large-scale problems.

Quasi-Newton Methods:

- Instead of calculating the Hessian matrix explicitly, quasi-Newton methods approximate the Hessian using only first-order information (gradients).
- **BFGS (Broyden-Fletcher-Goldfarb-Shanno)** and **L-BFGS (Limited-memory BFGS)** are popular quasi-Newton methods.
 - They approximate the Hessian based on gradient evaluations, making them much more efficient for large problems.

Truncated Newton (TNC):

- TNC uses a **conjugate gradient approach** to solve the optimization problem. It only computes an approximation of the Hessian and "truncates" the step when necessary to avoid going too far or using too much computational effort.
 - This makes it more scalable for large problems while maintaining the benefits of second-order methods.
-
- For machine learning applications, Newton's methods are often used in training models like **logistic regression**, where the optimization problem involves minimizing a log-likelihood function.
 - In these cases, algorithms like **BFGS** or **TNC** are used for efficient optimization especially in large-scale optimization.

Summary of Solvers in Scikit-learn Logistic Regression

Solver	Type	Key Characteristics	Use Case
<code>lbfgs</code>	Quasi-Newton	Efficient for large datasets, default solver	Binary and multinomial classification
<code>newton-cg</code>	Newton-based	High precision, better for multinomial problems	Multinomial logistic regression
<code>sag</code>	Gradient descent	Efficient for large datasets, supports L2 regularization	Large datasets with L2 regularization
<code>saga</code>	Gradient descent	Supports L1, L2, and elastic net regularization	Large datasets, sparse data, high-dimensional problems
<code>liblinear</code>	Coordinate descent	Good for small datasets, L1/L2 regularization	Binary classification, L1 regularization

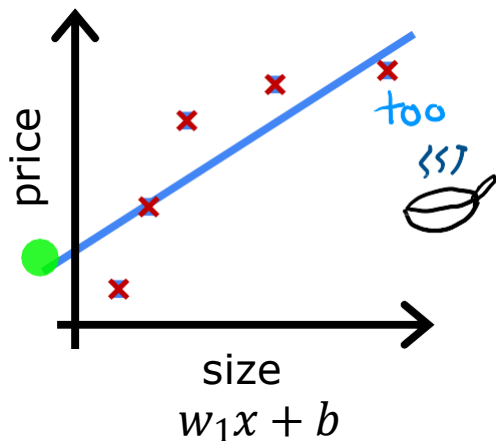


Regularization to Reduce Overfitting

The Problem of Overfitting



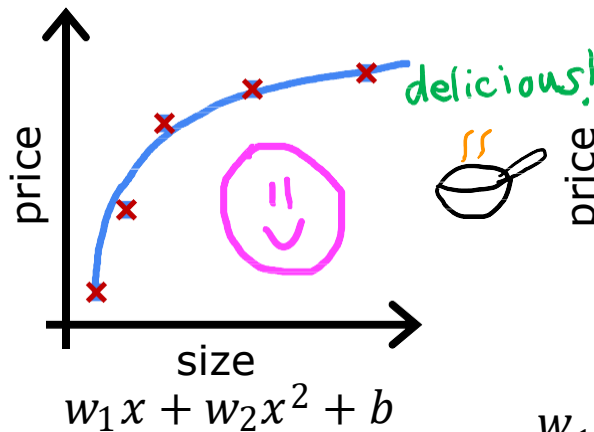
Regression example



Under fit

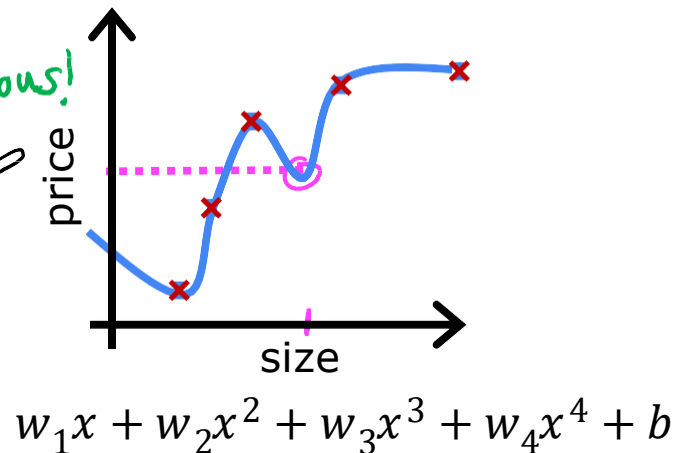
- Does not fit the training set well

high bias



- Fits training set pretty well

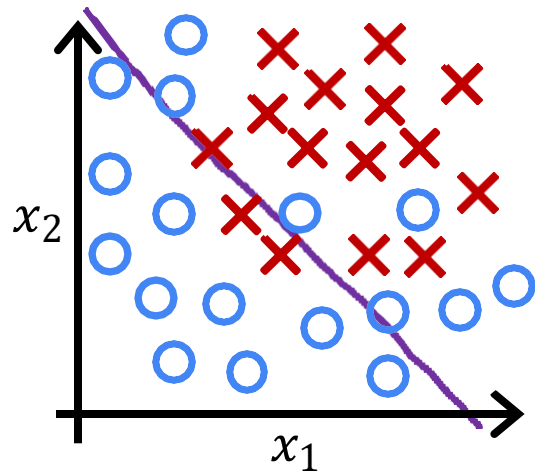
generalization



- Fits the training set extremely well

High variance

Classification

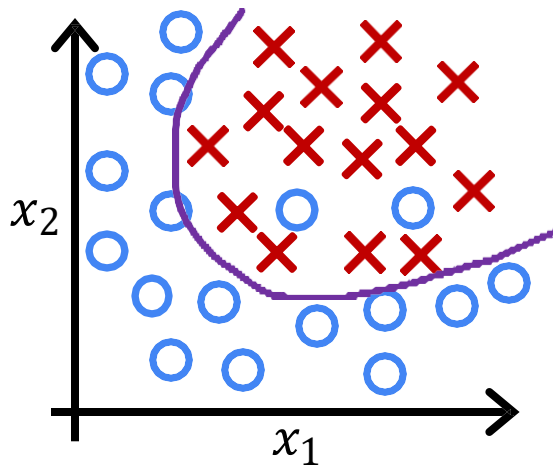


$$z = w_1x_1 + w_2x_2 + b$$

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

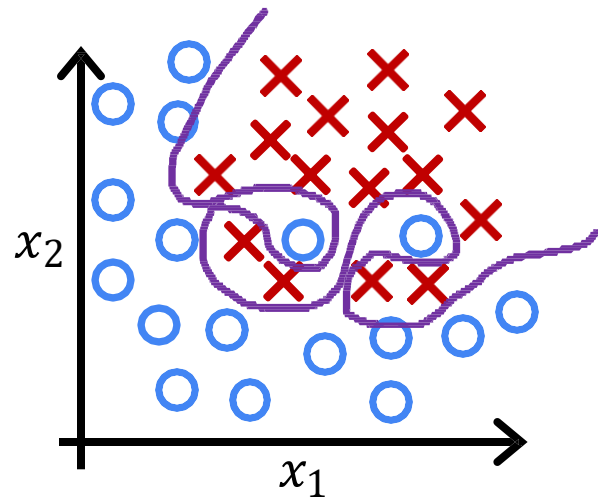
g is the sigmoid function

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



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

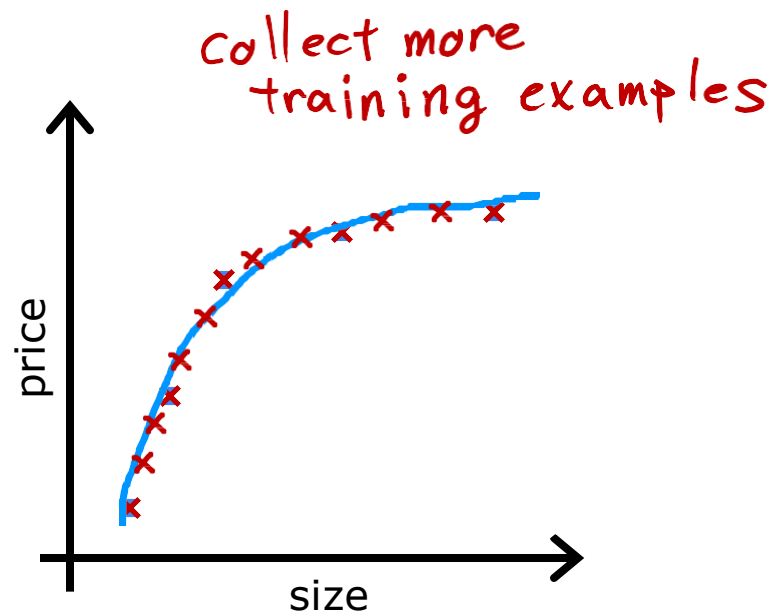
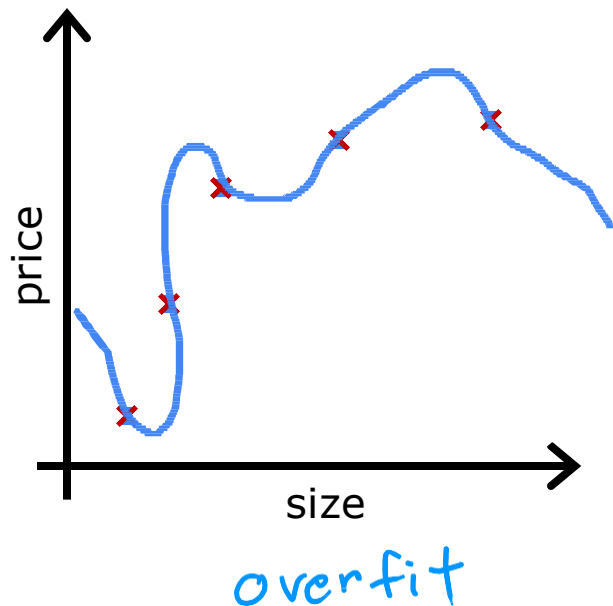
overfit



Regularization to Reduce Overfitting

Addressing Overfitting

Collect more training examples



Select features to include/exclude

size	bedrooms	floors	age	avg income	...	distance to coffee shop	price
x_1	x_2	x_3	x_4	x_5		x_{100}	y

all features



insufficient data



overfit

selected features

size

bedrooms

age

just right

feature selection

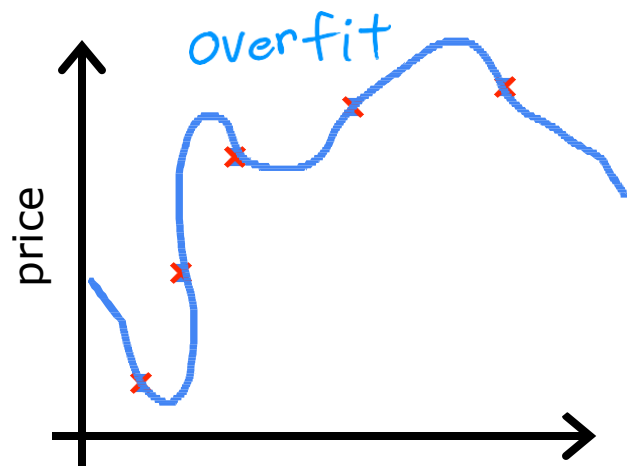
disadvantage



useful features
could be lost

Regularization

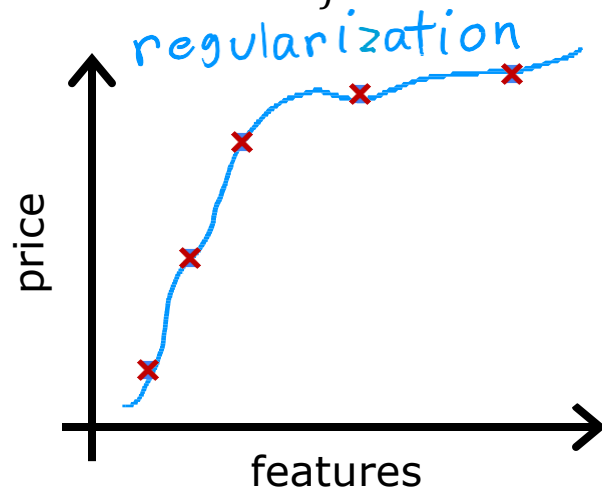
Reduce the size of parameters w_j



$$f(x) = 28x - 385x^2 + 39x^3 - 174x^4 + 100$$

large values w_j

eliminate feature



$$f(x) = 13x - 0.23x^2 + 0.000014x^3 - 0.0001x^4 + 10$$

small values for w_j

Addressing overfitting

Options

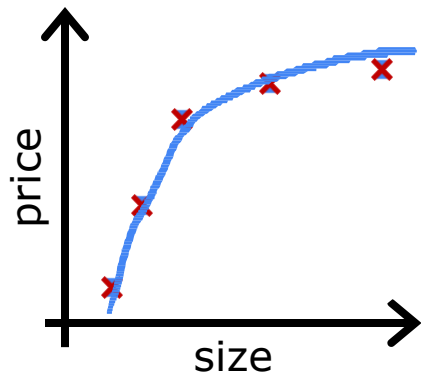
1. Collect more data
2. Select features
 - Feature selection
3. Reduce size of parameters
 - “Regularization”



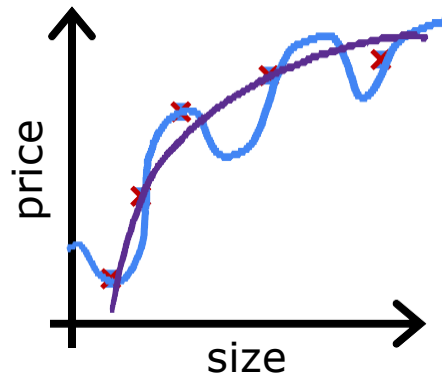
Regularization to Reduce Overfitting

Cost Function with
Regularization

Intuition



$$w_1x + w_2x^2 + b$$



$$w_1x + w_2x^2 + \underbrace{w_3x^3}_{\approx 0} + \underbrace{w_4x^4}_{\approx 0} + b$$

make w_3, w_4 really small (≈ 0)

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

$$+ 1000 \underbrace{w_3^2}_{0.001} + 1000 \underbrace{w_4^2}_{0.002}$$

Regularization

small values w_1, w_2, \dots, w_n, b

simpler model

less likely to overfit

$$w_3 \approx 0$$

$$w_4 \approx 0$$

size x_1	bedrooms x_2	floors x_3	age x_4	avg income x_5	...	distance to coffee shop x_{100}	price y
---------------	-------------------	-----------------	--------------	------------------------	-----	-----------------------------------------	--------------

$$w_1, w_1, w_2, \dots, w_{100}, b$$

n features

$$n = 100$$

$$J(\vec{w}, b) = \frac{1}{2m} \left[\sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{"lambda" regularization parameter}} + \underbrace{\frac{\lambda}{2m} b^2}_{\text{can include or exclude } b} \right]$$

regularization term

$\lambda > 0$

Regularization

$$\min_{\vec{w}, b} J(\vec{w}, b) = \min_{\vec{w}, b} \left[\underbrace{\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2}_{\text{mean squared error}} + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{regularization term}} \right]$$

fit data

Keep w_j small

λ balances both goals

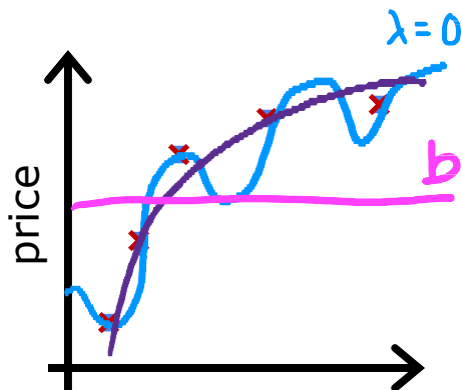
choose $\lambda = 10^{10}$

$$f_{\vec{w}, b}(\vec{x}) = \cancel{w_1}x + \cancel{w_2}x^2 + \cancel{w_3}x^3 + \cancel{w_4}x^4 + b$$

≈ 0 ≈ 0 ≈ 0 ≈ 0

$$f(x) = b$$

choose λ





Regularization to Reduce Overfitting

Regularized Linear
Regression

Regularized linear regression

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

Gradient descent

repeat {

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

$j = 1, \dots, n$

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

} simultaneous update

Implementing gradient descent

repeat {

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

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

} simultaneous update $j = 1 \dots n$

Implementing gradient descent

repeat {

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

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

} simultaneous update $j = 1 \dots n$

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

shrink w_j

$$\alpha \frac{\lambda}{m} = 0.01 \frac{1}{50} = 0.0002$$

$$w_j (1 - 0.0002) = 0.9998 w_j$$

How we get the derivative term (optional)

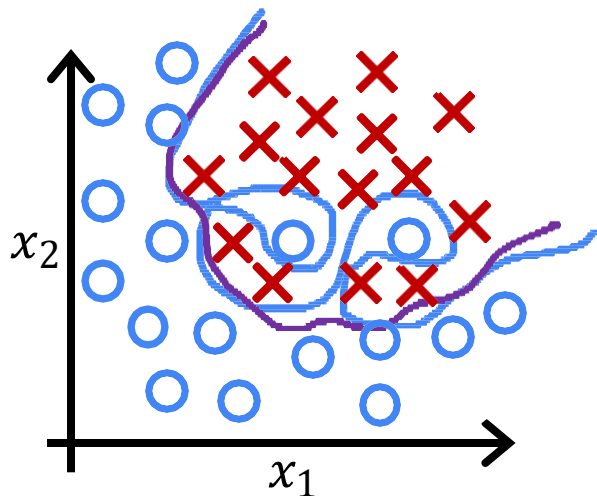
$$\begin{aligned}
 \frac{\partial}{\partial w_j} J(\vec{w}, b) &= \frac{d}{dw_j} \left[\frac{1}{2m} \sum_{i=1}^m \left(\underbrace{f(\vec{x}^{(i)})}_{\vec{w} \cdot \vec{x}^{(i)} + b} - y^{(i)} \right)^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right] \\
 &= \cancel{\frac{1}{2m} \sum_{i=1}^m \left[(\vec{w}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



Regularization to Reduce Overfitting

Regularized Logistic
Regression

Regularized logistic regression



$$z = w_1x_1 + w_2x_2 + w_3x_1^2x_2 + w_4x_1^2x_2^2 + w_5x_1^2x_2^3 + \dots + b$$

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

Cost function

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

$\min_{\vec{w}, b} J(\vec{w}, b) \rightarrow w_j \downarrow$

Regularized logistic regression

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

$\min_{\vec{w}, b}$

Gradient descent

repeat {

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

$j=1 \dots n$

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

}

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

logistic regression

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

don't have to regularize b