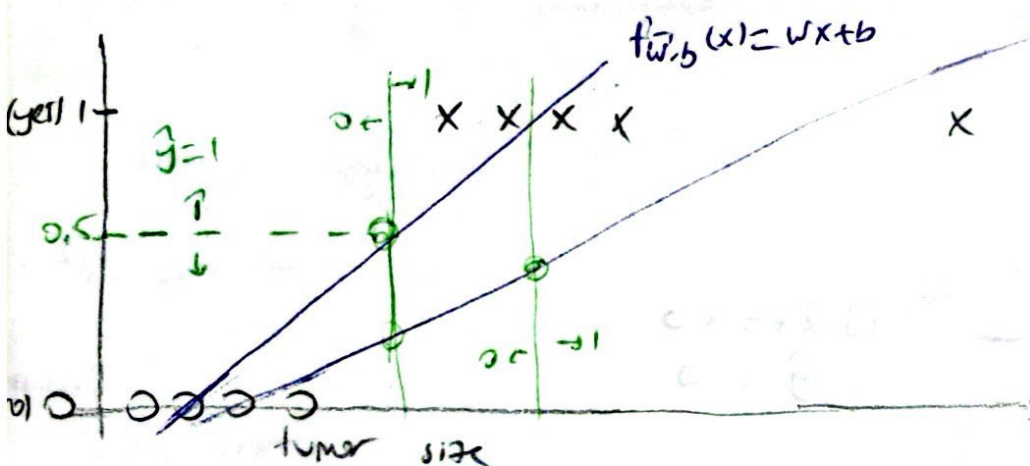*Linear regression is not a good idea algorithm for classification problems. → logistic regression

## Classification

• Is this email spam?
• Is the tumor malignant

Answer "y"

no     yes
(0)    (1)

→ y can only be one of two values. → binary classification



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

• if $f_{w,b}(x) < 0.5 → \hat{y} = 0$
• if $f_{w,b}(x) \geqslant 0.5 → \hat{y} = 1$

→ linear regression is not god idea fo this classification

# Logistic Regression

Want outputs between 0 and 1



- Sigmoid function
- logistic function

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

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

- "probability" that class is 1

$x$ : tumor size

$y$ : 0 (not malignant)
     1 (malignant)

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

70% chance that $y$ is 1
30% chance that $y$ is 0

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

    yes: $\hat{y} = 1$,    no: $\hat{y} = 0$

* When is $f_{\vec{w},b}(\vec{x}) \geq 0.5$
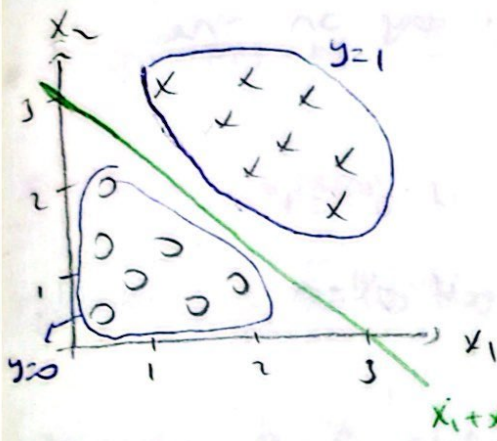
      $g(z) \geq 0.5$

       $z \geq 0$

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

      $\hat{y} = 1$                $\hat{y} = 0$

# Decision boundary:



$$f_{\vec{w},b}(\vec{x}) = g(z) = g(w_1 x_1 + w_2 x_2 + b)$$
$$\downarrow \quad \downarrow \quad \downarrow$$
$$1 \quad 1 \quad -3$$

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

$$\boxed{x_1 + x_2 = 3} \rightarrow \text{decision boundary}$$

**\*** Logistic regression con learn to fit pretty complex data.

**\*** The decision boundary for logistic regression will always be linear, will always be a straight line.

## Cost Func for Logistic Regression

$$f_{\vec{w},b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}} \qquad \text{target } y \text{ is } 0 \text{ or } 1$$
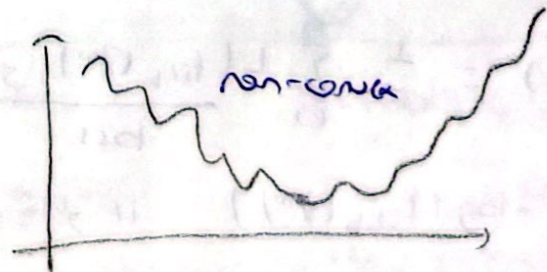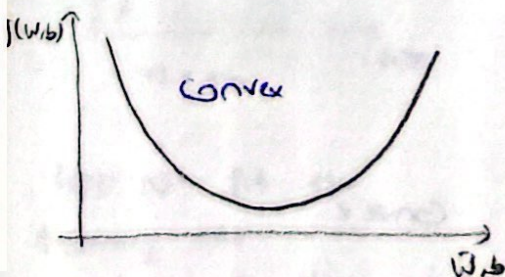
$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^{m} \frac{1}{2} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right)^2$$

### linear regression

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



convex

### logistic regression

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



non-convex

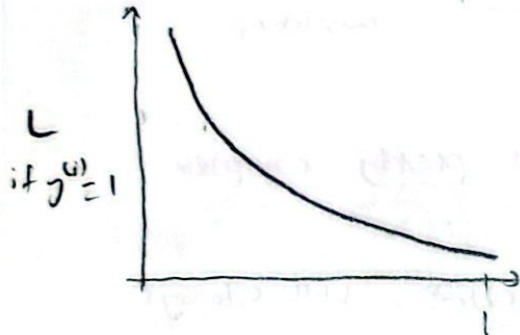**\*** This squared error cost fnc is not a good choice for logistic regression → there are lots of local minima, non-convex.

There will be different cost fnc.

# Logistic loss function:

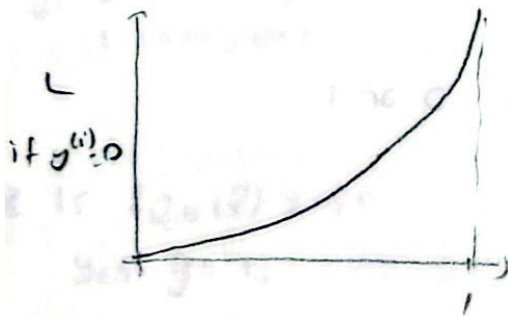The loss function measures how well you're doing on one training examples

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



- $f_{\vec{w},b}(x^{(i)}) \to 1$ then loss $\to 0$
- $f_{\vec{w},b}(\vec{x}^{(i)}) \to 0$ then loss $\to \infty$

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



- $f_{\vec{w},b}(\vec{x}^{(i)}) \to 0$ then loss $0$
- $f_{\vec{w},b}(\vec{x}^{(i)}) \to 1$ then loss $\infty$

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

## Cost

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

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

$\implies$ Convex

- can reach a global minimum

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

y can only be 0 or 1

## Gradient descent

$$J(\vec{w},b) = -\frac{1}{M} \sum_{i=1}^{M} \left[ y^{(i)} \log(f_{\vec{w},b}(\vec{x}^{(i)})) + (1-y^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}^{(i)})) \right]$$
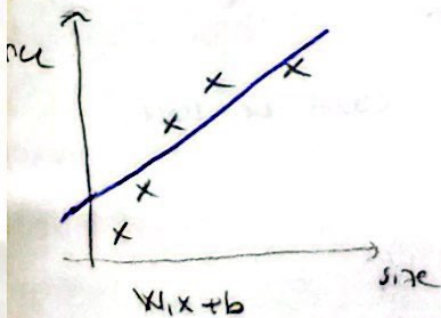
• $w_j = w_T - a \cdot \dfrac{\partial}{\partial w_T} J(\vec{w},b)$

$\longrightarrow \dfrac{1}{M} \sum_{i=1}^{M} (f_{w,b}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$    } simultaneously update

• $b = b - a \cdot \dfrac{\partial}{\partial b} J(\vec{w},b)$

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

⟹ Some concepts:
• monitor gradient descent (learning curve)
• Vectorized implementation.
• Feature scaling

## The Problem of Overfitting

Regression



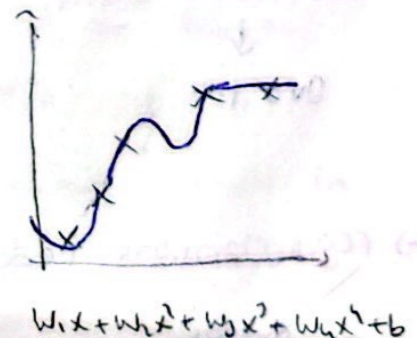| Wix +b | Wix + W₂x² + b | W₁x + W₂x² + W₃x³ + W₄x⁴ + b |
|---|---|---|

• Does not fit the training set

− Underfit

• high bias

• Fits training set pretty well
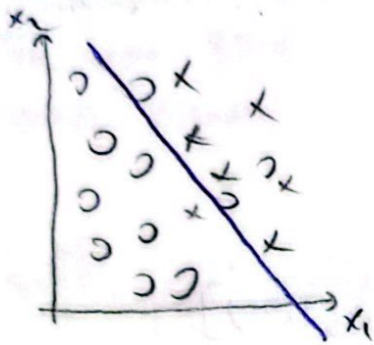
Just right

generalization

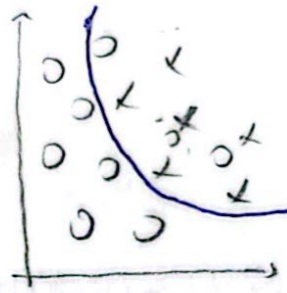• Fits the training set extremely well

− overfit

high variance

• machine learning goal → a model that has neither high bias nor high variance
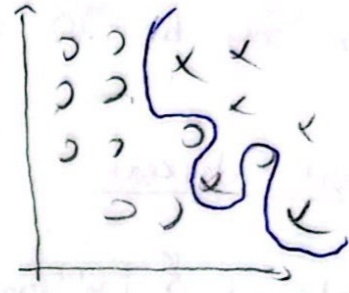
# Classification



- underfit          - Just right          - overfit

- Its model despite doing very well on the traning set, doesn't look like it'll generalize well to new examples.

## Addressing overfitting

If there is overfit prediction model

→ collect more traning data  ,  learning algorithm will learn to fit a function that is less wiggly.

* more data isn't always a option.

| all features | selected features |
|---|---|
| + | ↓ |
| insufficient data | JUST right — |
| ↓ | feature selection |
| overfit | * Useful features could be lost |

→ regularizations reduce the size of parameters $w_j$



$f(x) = 28x + 39x^3 + 109 ...$

$f(x) = 13x + 0.000014x^3 - 0.001x^4 + 10$

* regularization:· keep all of your features
                    , but they just prevents the features from having
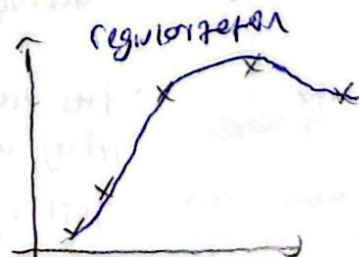                    on overly large effect.
                    · Sometimes can, cause overfitting


Addressing overfitting
Options 1) collect more data

        2) Select features
            - feature selection

        3) Reduce size of parameters - "Regularization"


Cost function with regularization

• $w_1 x + w_2 x^2 + b$                    • $w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + y$

make $w_3, w_4$ very small ($\approx 0$)

$$\min_{\vec{w}, b} \frac{1}{2m} \sum_{i=1}^{m} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + 1000 \, w_3^2 + 1000 \, w_4^2$$
                                            \ 0,001          0,001

* Simpler model → less likely to overfit

& regularization implemented → to penalize all of the features
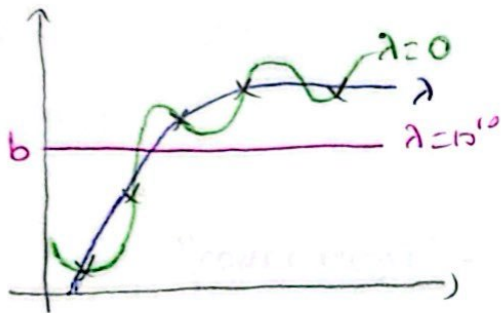                                        or more precisely
(
& and it's possible to show that this will usually result in
fitting a smoother simpler, less weckly function that's less
prone to overfitting.


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

$\min J(\vec{w}, b)$ → The algorithm also tries to keep $w_J$ small,
                        which will tend to reduce overfitting.

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

$\underbrace{\hphantom{fit data xxxxxx}}_{\text{fit data}}$  $\underbrace{\hphantom{keep w small}}_{\text{keep } w_j \text{ small}}$

$\lambda$ balance both pals



$\lambda: $ balanced

- minimizing the mean squared error
- keeping parameters small

## Regularized linear regression

implementing gradient descent

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

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

## Regularized logistic regression

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

$$w_j = w_j - a \cdot \frac{\partial}{\partial w_j} J(w,b)$$

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

$\Big\}$ not same linear