# Linear regression

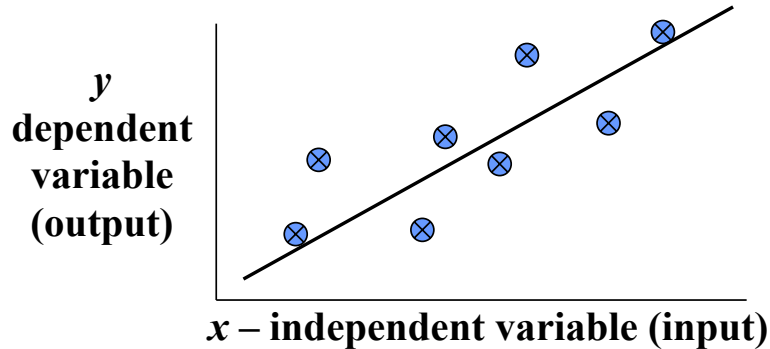## Hung-Hsuan Chen

# Parametric supervised learning

- Notation
  - Features    $x$
  - Targets     $y$
  - Predictions $\hat{y}$
  - Parameters $\theta$  (to be learned)

**Learning algorithm**

Change $\theta$
Improve performance

**Program  ("Learner")**

Characterized by
some "parameters"   $\theta$

Procedure (using $\theta$)
that outputs a prediction

**Training data
(examples)**

**Features**

**Feedback /
Target values**

**Score performance
("cost function")**

# Linear regression



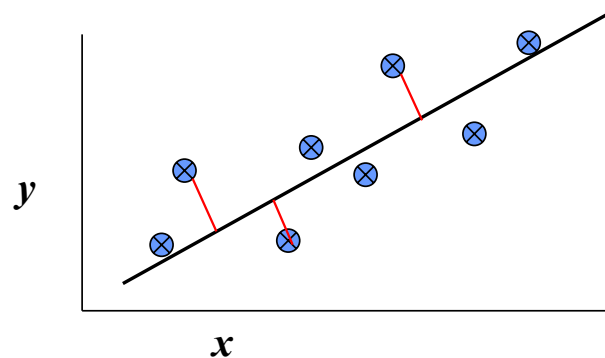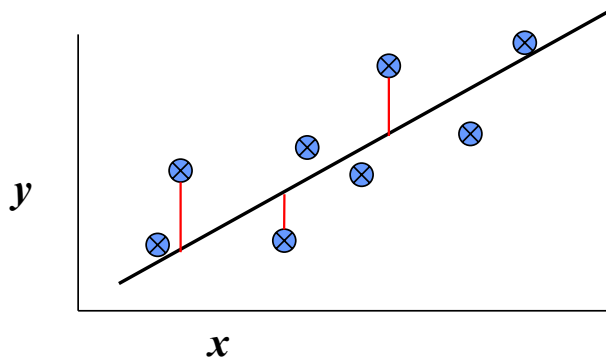$y$ dependent variable (output)
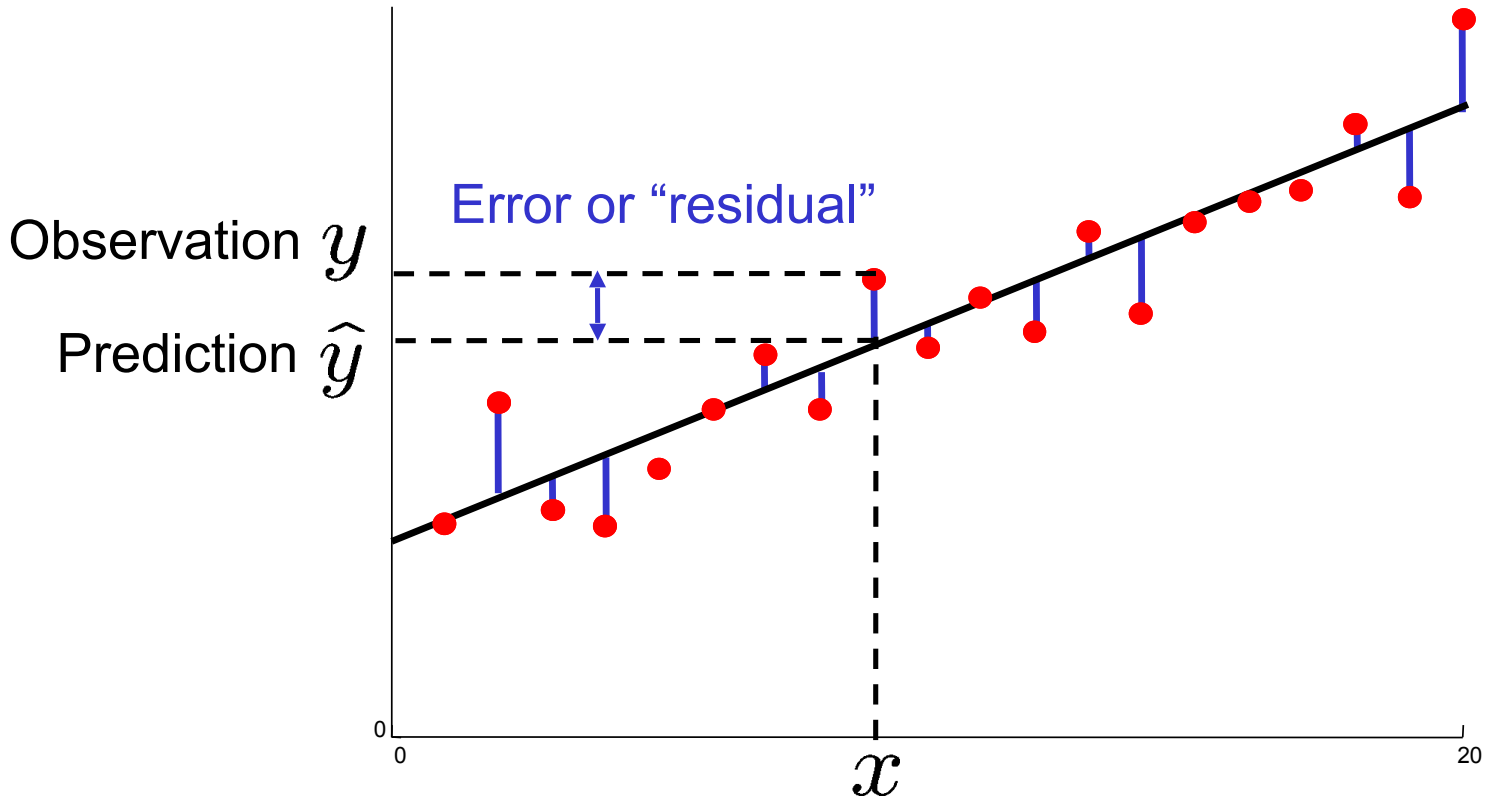
$x$ – independent variable (input)

"Predictor":

$$\hat{y} = \theta_0 + \theta_1 x$$

- Define the form of the function f(x) explicitly
  - i.e., $\hat{y} = \theta_0 + \theta_1 x$ in this case
  - If you believe that *x* and *y* have a **non-linear** relationship, you should assume a non-linear f(x)
- Find a good f(x) within that family
  - i.e. find good $\theta_0$ and $\theta_1$ such that $\hat{y} \approx y$ in this case

# Quiz: which one should be an error measurement?

# Measuring error

# Simple linear regression

- For now, assume just one (input) independent variable *x*, and one (output) dependent variable *y*
  - Multiple linear regression assumes an input vector **x**
- We will "fit" the points with a line
  - In a high dimensional space, we will fit the points with a hyper-plane
- Which line should we use?
  - Choose an objective function
  - Typically, we choose sum squared error (SSE) (or its variants)

$$\frac{1}{n} \sum \left( \hat{y}_i - y_i \right)^2 = \frac{1}{n} \sum \left( \text{residual}_i \right)^2$$

  - Choices with the same result: $\frac{1}{2n} \sum \left( \hat{y}_i - y_i \right)^2$ or $\frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2$

  - You may choose other objective functions
    - More to come in future lectures

why $\frac{1}{2}$？微分後 好看

# How to "learn" the parameters?

- For the 2-$d$ problem there are coefficients for the bias and the independent variable ($y$-intercept and slope)
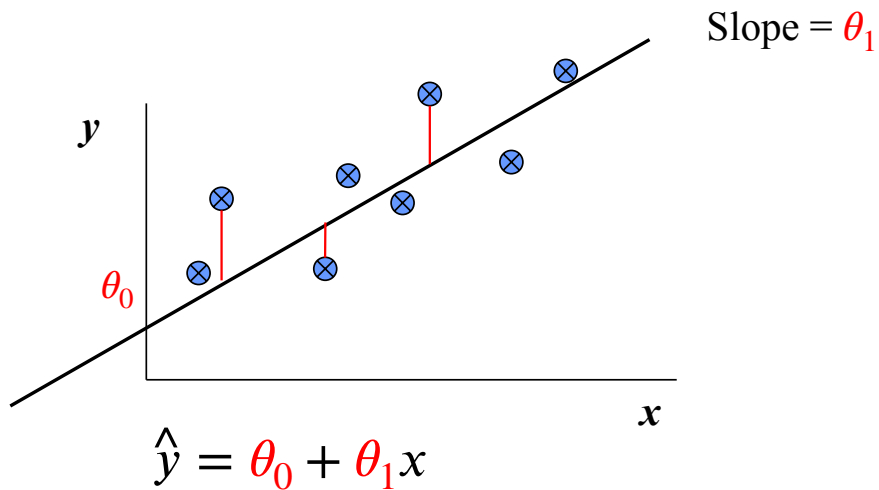  - $\hat{y} = \theta_0 + \theta_1 x$

- The value of the coefficients which minimize the objective function:

  - $$\theta_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - \left( \sum x_i \right)^2}$$

  - $$\theta_0 = \frac{\sum y_i - \theta_1 \sum x_i}{n}$$

  - $n$: number of training instances

# Visualizing $\theta_0$ and $\theta_1$

Slope = $\theta_1$

$y$

$\theta_0$

$x$

$$\hat{y} = \theta_0 + \theta_1 x$$

# How to derive the value of the coefficients?

- $J(\theta_0, \theta_1) = \dfrac{1}{2n} \sum \left( \hat{y}_i - y_i \right)^2$

$$= \dfrac{1}{2n} \sum \left( (\theta_0 + \theta_1 x_i) - y_i \right)^2 ,$$
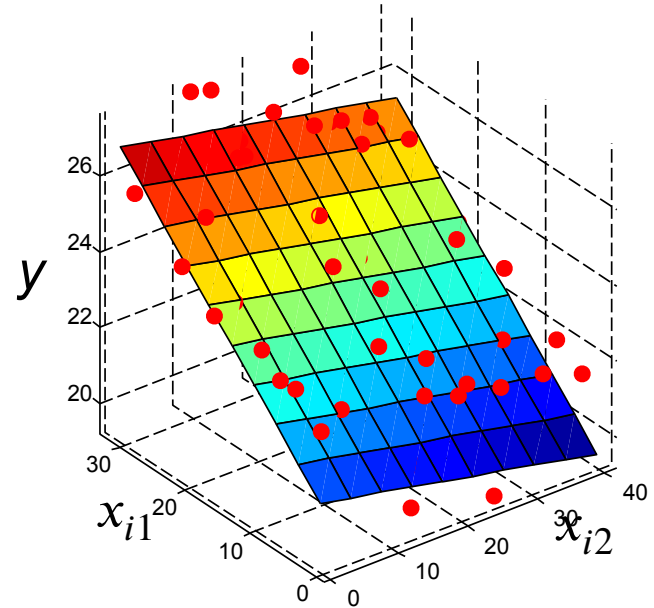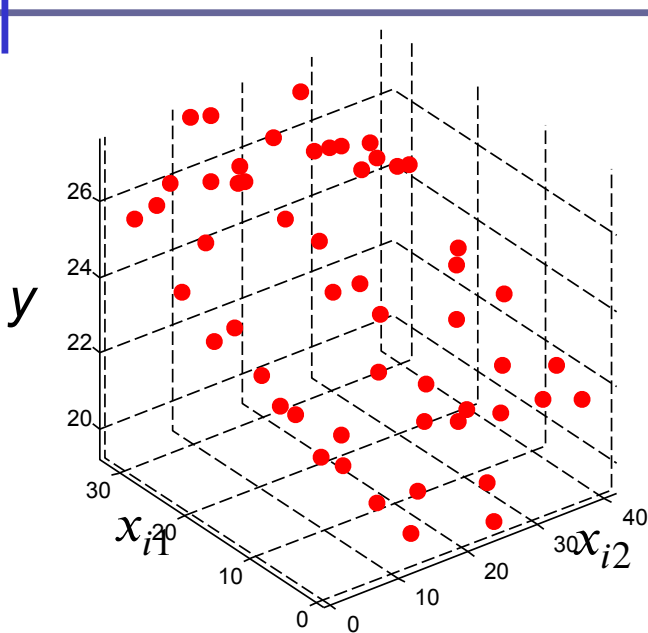
- Set

  - $\dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = 0$ and

  - $\dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = 0$

  to solve $\theta_0$ and $\theta_1$

# More dimensions?



$$\hat{y}(\mathbf{x}_i) = \theta_0 + \theta_1 x_{i,1} + \theta_2 x_{i,2} + \ldots + \theta_d x_{i,d} = \boldsymbol{\theta}^T \mathbf{x}_i, \text{ where}$$

$$\mathbf{x}_i = \left[ 1, x_{i,1}, x_{i,2}, \ldots, x_{i,d} \right]^T$$

$$\boldsymbol{\theta} = \left[ \theta_0, \theta_1, \theta_2, \ldots, \theta_d \right]^T$$

$d$: # of features, $x_{i,j}$: the $i$th training instance's $j$th feature

# Multiple linear regression

- $n$: the number of training instances
- $d$: the number of features
- Training instances:

  $$X = \begin{bmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{bmatrix}, \; y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

  – (We assume no coefficient parameter here)

- Find $\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$ such that $\left( \hat{y} - y \right)^T \left( \hat{y} - y \right)$ is minimized,

  where
  – $\hat{y} = X\theta$

- The solution is $\theta = \left( X^T X \right)^{-1} X^T y$

# How to derive the value of the coefficient vector?

- $J(\theta) = \dfrac{1}{2}\left(\hat{y} - y\right)^T \left(\hat{y} - y\right)$

  $= \dfrac{1}{2}\left(X\theta - y\right)^T \left(X\theta - y\right)$

- Set
  - $\nabla J(\theta) = 0$

  to solve $\theta$

# Another way to find $\theta$ for the linear regression problem

- Goal: find $\theta_0$ and $\theta_1$ to minimize

  $$- J(\theta_0, \theta_1) = \frac{1}{2n} \sum \left( \hat{y}_i - y_i \right)^2$$

  $$= \frac{1}{2n} \sum \left( (\theta_0 + \theta_1 x_i) - y_i \right)^2$$

- Procedure
  1. Start with $\theta_0 = r_0$ (a random number) and $\theta_1 = r_1$ (another random number)
  2. Slightly move $\theta_0$ and $\theta_1$ to reduce $J(\theta_0, \theta_1)$
  3. Keep doing step 2 until converged
- Question: how to move $\theta_0$ and $\theta_1$?

# Plotting the objective function (3D and contour plots)



Keep moving "downward" to reach the minimum
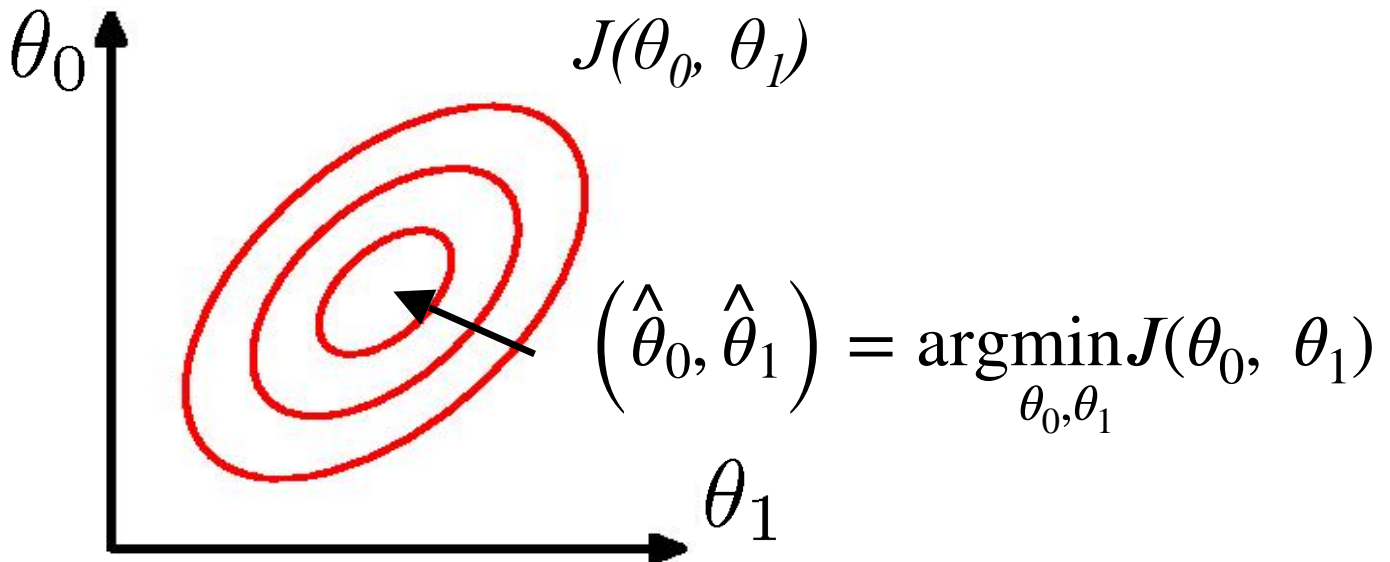
# Finding good parameters

- Want to find parameters which minimize our error…

- Think of a cost "surface": error residual for that θ…



$$J(\theta_0, \theta_1)$$

$$\left(\hat{\theta}_0, \hat{\theta}_1\right) = \operatorname*{argmin}_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

# Gradient descent

$\theta_0^{(k)} = 5$

$\theta_1^{(k)} = -4$

$\hat{y} = \theta_0 + \theta_1 x$

$J = \left( y - \left( \theta_0 + \theta_1 x \right) \right)^2$

$\dfrac{\partial J}{\partial \theta_0} = 2 \left( \quad \right) (-1)$

$\dfrac{\partial J}{\partial \theta_1} = 2 \left( \quad \right) (-x)$

- Procedure
  1. Start with random values
     - $\boldsymbol{\theta} = \boldsymbol{\theta}^{(0)} = (\theta_0^{(0)}, \theta_1^{(0)}, \ldots, \theta_d^{(0)})$
  2. Slightly move $\theta_0, \ldots \theta_d$ to reduce $J(\boldsymbol{\theta})$
     - $$\theta_i^{(k+1)} = \theta_i^{(k)} - \alpha \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_i} \bigg|_{\boldsymbol{\theta} = \boldsymbol{\theta}^{(k)}}$$

       α is a small positive number

     - $k = k + 1$
  3. Keep doing step 2 until con~~verges~~

$$\frac{\partial}{\partial \theta_i} J(\theta) = \frac{1}{n} \sum_j \left( \hat{y}_j - y_j \right) x_{ji},$$
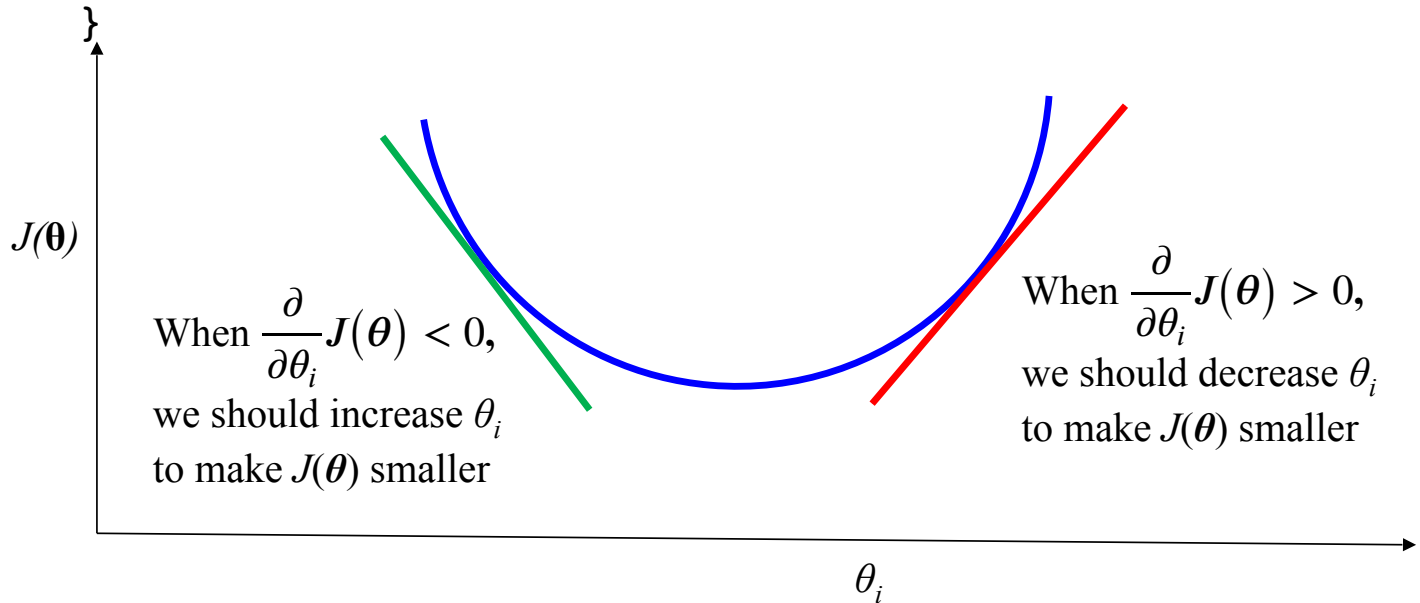
$\dfrac{1}{n}$ is usually ignored

# Why gradient descent work?

```
Repeat until converge {
   for (i=1,2, ..., d){
```

$$\theta_i := \theta_i - \alpha \frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) = \theta_i - \alpha \sum_j \left(\hat{y}_j - y_j\right) x_{ji}$$

```
   }
}
```

$J(\boldsymbol{\theta})$

When $\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) < 0,$
we should increase $\theta_i$
to make $J(\boldsymbol{\theta})$ smaller

When $\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) > 0,$
we should decrease $\theta_i$
to make $J(\boldsymbol{\theta})$ smaller

$\theta_i$

# Gradient descent (one feature and two features)

# Overly large learning rate may not lead to converge

- $$\theta_i^{(k+1)} := \theta_i^{(k)} - \alpha \frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}^{(k)})$$

  - $\alpha$: learning rate
  - Often $\alpha \in [0.001, \ 1]$
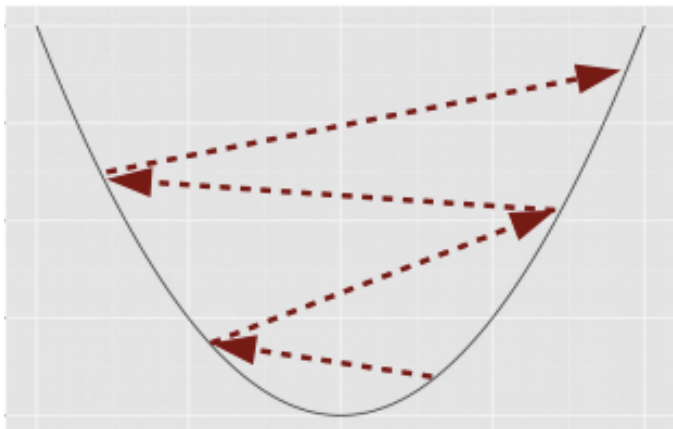  - Shrink $\alpha$ as $k$ becomes larger

# We now have two tools to solve linear regression problem

- **Close form solution**

    $$\theta = \left(X^T X\right)^{-1} X^T y$$

- Gradient descent

```
Repeat until converge {
  for (i=1,2, ..., d){
```
$$\theta_i = \theta_i - \alpha \sum_j \left(\hat{y}_j - y_j\right) x_{ji}$$
```
  }
}
```

# Which one should we use?

# Consider using close form solution

- $\theta = \left( X^T X \right)^{-1} X^T y$
- Required space:
  - Even if X is sparse, $(X^TX)^{-1}$ may not be sparse
  - If n=1M, d=100K, $(X^TX)^{-1}$'s dimension is $(100k)^2 = 10^{10}$
    - If 8 bytes per entry, storing $(X^TX)^{-1}$ requires 80GB ➜ typically infeasible on a single machine
- Computation time:
  - Computing matrix multiplication $X^TX$ takes $O(nd^2)$
  - Let $\left( X^T X \right) = P$, computing $P^{-1}$ takes $O(d^{2.373})$ to $O(d^3)$
    - Depending on the inversion algorithm
  - Let $(X^TX)^{-1} = Q$, computing $QX^T$ takes $O(nd^2)$
  - Let $\left( X^T X \right)^{-1} X^T = R$, computing Ry takes $O(nd)$
  - Total: $O(nd^2+d^{2.373}+nd^2)$
- If d is small, close form solution is probably acceptable

n: number of
Instance

d: number of
Features

# Consider using Gradient Descent

```
Repeat until converge {
  for (j=1, 2, …, n){
```
$$z_j = \left( \boldsymbol{\theta}^T \mathbf{x}_j - y_j \right)$$
```
  }
  for (i=1, 2, ..., d){
```
$$\theta_i = \theta_i - \alpha \sum_j z_j x_{ji}$$
```
  }
}
```

- Required space:
  - $\mathbf{X} = [x_{ij}]$, $\mathbf{y}$, $\boldsymbol{\theta}$, and $\mathbf{z}$
  - If n=1M, d=100K, and X is sparse
    - If 8 bytes per entry, we need: (nnz(X) +1M+100K+1M)*8bytes
    - Much more efficient than 80GB
- Computation time:

$$- \quad z_j = \left( \boldsymbol{\theta}^T \mathbf{x}_j - y_j \right) \text{ for all } j:$$

  O($nd$)

$$- \quad \theta_i := \theta_i - \alpha \sum_j z_j x_{ji} \text{ for all } i:$$

  O($nd$)
  - Outer loop: Assume $T$ iterations
  - Total: O($Tnd$)
    - Usually more efficient than O($nd^2 + d^{2.373} + nd^2$)

# Close form vs gradient descent

- If the number of features is small, close form solution is probably acceptable

- However, if the number of features is large, using gradient descent is more efficient in both memory usage and computing time

- Moreover, gradient descent is capable of solving more complex optimization problem
  - In many cases, $\dfrac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{0}$ has no closed–form solution
    - But we can still apply gradient descent ☺
  - More to come in future lectures

# Quiz

- If the features are all (categorical)
  - Can you apply linear regression?
  - Can you apply decision tree?
  - Can you apply knn?
- If all features are numerical, which one runs faster during "test" phase?
  - Linear regression? Decision tree? KNN?

Lg    little faster than   Dt

# When to stop updating?

- Possible stopping criteria:
    - Improvement drops (e.g., < 0)
    - Reached small error
    - Achieved predefined # of iterations
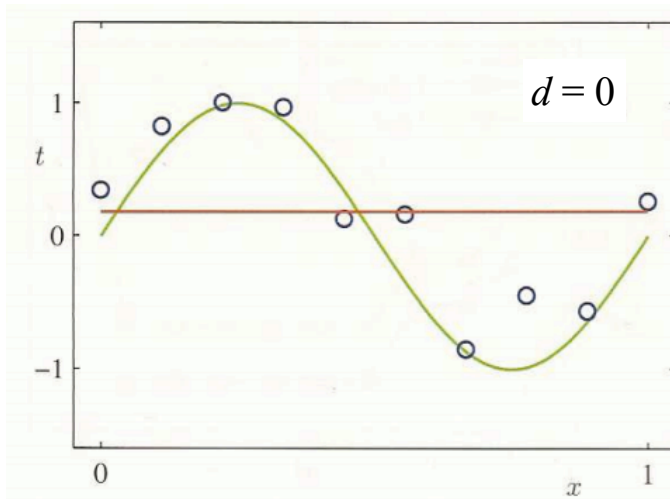    - No time to train anymore

# Fitting non-linear data to linear model

- $\hat{y}(\mathbf{x_i}) = \theta_0 + \theta_1 x_{i,1} + \theta_2 x_{i,2} + \dots + \theta_d x_{i,d}$

  - Linear model

- We may generate the higher degree terms as the new features

  - $\hat{y}(\mathbf{x_i}) = (\theta_0 + \theta_1 x_{i,1} + \theta_2 x_{i,2} + \dots + \theta_d x_{i,d}) +$
    $(\theta_{1,2} x_{i,1} x_{i,2} + \theta_{1,3} x_{i,1} x_{i,3} + \dots + \theta_{d-1,d}\ x_{i,d-1} x_{i,d}) +$
    $(\theta_{1,1} x_{i,1}^2 + \theta_{2,2} x_{i,2}^2 + \dots + \theta_{d,d} x_{i,d}^2)$
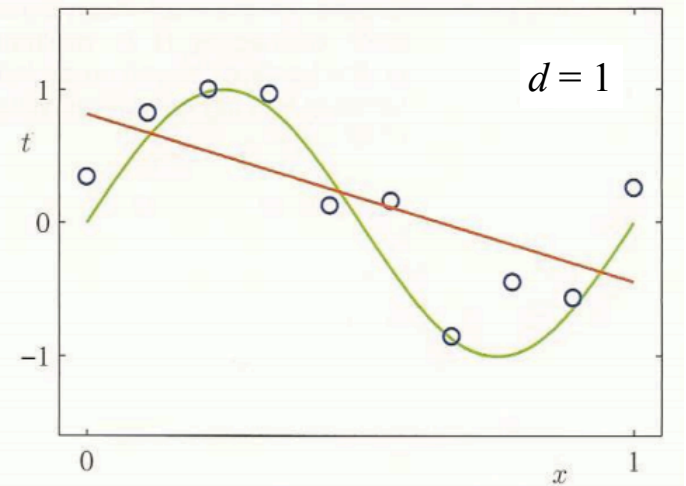
# Overfitting

- Overfitting occurs when a model captures idiosyncrasies of the input data, rather than generalizing
  - Too many parameters relative to the amount of training data
- For example, an order-$N$ polynomial can perfectly fit to $N+1$ data points

# Target: $\sin(2\pi x)$ + noise



$d = 0$

$d = 1$
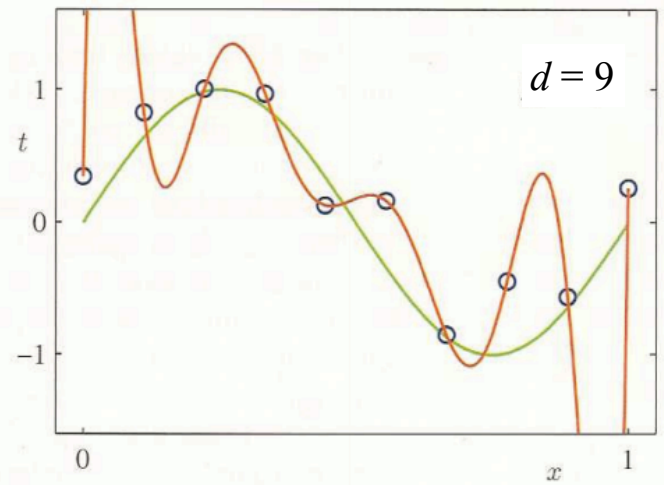
$\hat{y}_i = \theta_0$

$\hat{y}_i = \theta_0 + \theta_1 x$

# Target: $\sin(2\pi x) +$ noise (overfitting)

$$\hat{y_i} = \theta_0 + \sum_{i=1}^{9} \theta_i x^i$$



$d = 3$

$d = 9$

$$\hat{y_i} = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

# Observation

- In the linear regression model, overfitting is characterized by large parameters

|  | $d = 0$ | $d = 1$ | $d = 3$ | $d = 9$ |
|---|---|---|---|---|
|  | 0.19 | 0.82 | 0.31 | 0.35 |
|  |  | -1.27 | 7.99 | 232.37 |
|  |  |  | -25.43 | -5321.83 |
|  |  |  | 17.37 | 48568.31 |
|  |  |  |  | 231639.30 |
|  |  |  |  | 640042.26 |
|  |  |  |  | -1061800.52 |
|  |  |  |  | 1042400.18 |
|  |  |  |  | -557682.99 |
|  |  |  |  | 125201.43 |

# Regularization to avoid overfitting

- Introduce a penalty term for the size of the weights
- Un-regularized regression (the original objective function)

  - $$J(\mathbf{\theta}) = \frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 = \frac{1}{2} \sum \left( \mathbf{\theta}^T \mathbf{x}_i - y_i \right)^2$$

  - $$\mathbf{\theta} := \operatorname*{argmin}_{\mathbf{\theta}} \frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2$$

- Regularized regression (enforce the solution to have low L2-norm of $\mathbf{\theta}$)

  - $$\mathbf{\theta} := \operatorname*{argmin}_{\mathbf{\theta}} \frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 \text{ such that } \left\| \mathbf{\theta} \right\|^2 \leq K$$

$$\theta_1^2 + \theta_2^2 \cdots \theta_n^2$$

# Regularization to avoid overfitting

- Original objective function: minimizing the training error

$$J(\boldsymbol{\theta}) = \frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 = \frac{1}{2} \sum \left( \boldsymbol{\theta}^T \mathbf{x}_i - y_i \right)^2$$

- New target

$$\boldsymbol{\theta} := \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 \text{ such that } \|\boldsymbol{\theta}\|^2 \leq K$$

- This is equivalent to the following problem with some $\lambda$

$$\boldsymbol{\theta} := \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left[ \frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 \right] \rightarrow \text{hyperporaneter}$$

- New objective function: minimizing training error and the L2-norm of $\boldsymbol{\theta}$ simultaneously

  - Solution: $\boldsymbol{\theta} = \left( \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$

＊ 沒λ大
得到的 θ 小

# Regularized linear regression

- Regularized linear regression

  - $\boldsymbol{\theta} := \underset{\boldsymbol{\theta}}{\arg\min} \left[ \dfrac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 + R(\boldsymbol{\theta}) \right]$

  - $R(\boldsymbol{\theta})$: regularization

- L2-regularization (**Ridge regression**)

  - $\boldsymbol{\theta} := \underset{\boldsymbol{\theta}}{\arg\min} \left[ \dfrac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 + \lambda \|\boldsymbol{\theta}\|^2 \right]$

    $\rightarrow \theta_1^2 + \theta_2^2 \dots$

- L1-regularization (**Lasso**)

  - $\boldsymbol{\theta} := \underset{\boldsymbol{\theta}}{\arg\min} \left[ \dfrac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 + \lambda \|\boldsymbol{\theta}\|_1 \right]$

    $\rightarrow |\theta_1| + |\theta_2| + \dots$

# Combining Ridge and Lasso

- Elastic net regularization

  $$\boldsymbol{\theta} := \operatorname*{argmin}_{\boldsymbol{\theta}} \left[ \sum \left( \hat{y}_i - y_i \right)^2 + \lambda_1 \|\boldsymbol{\theta}\|_1 + \lambda_2 \|\boldsymbol{\theta}\|^2 \right]$$

  - When $\lambda_1 = 0$ and $\lambda_2 > 0$ ➔  Ridge regression
  - When $\lambda_1 > 0$ and $\lambda_2 = 0$ ➔  Lasso

# How to select $\lambda$?

- Split the data into training and test datasets
  - Based on the training data, train the models by different $\lambda$'s
  - Based on the test data, calculate the test performance of all the models (of different $\lambda$'s)
  - Select the $\lambda$ with the best test performance
- Cross validation
  - More to come in future lectures

# Properties of Ridge and Lasso

- Ridge
  每り feature 皆有 effect
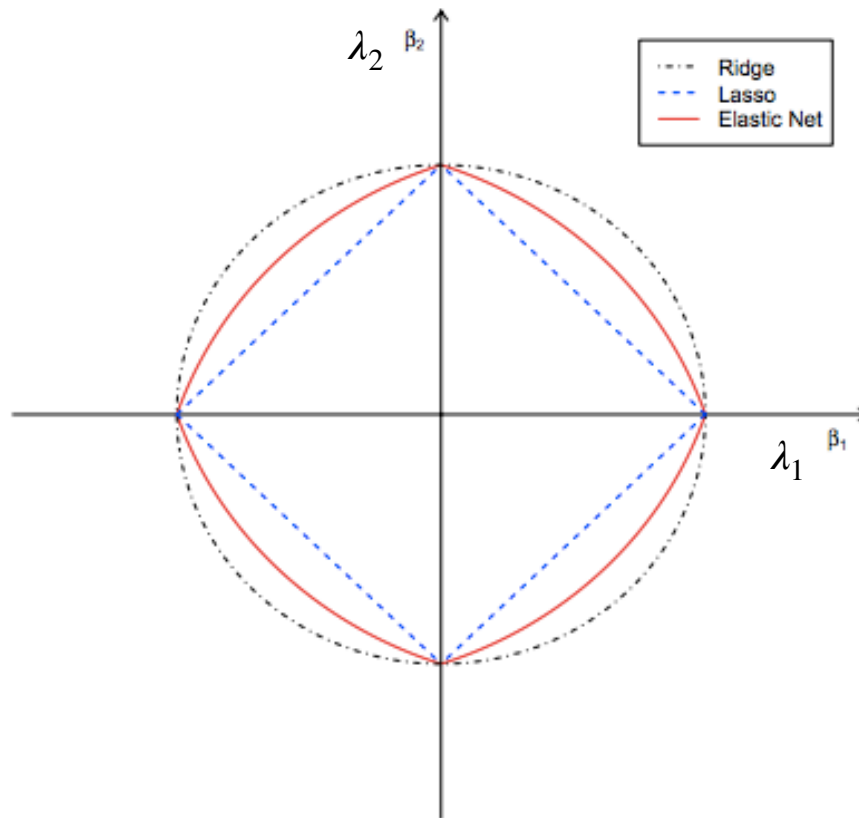  - Good if many features have small/medium sized effects

- Lasso
  只有 特定 feature 有 effect
  - Good if only a few features with a medium/ large effect
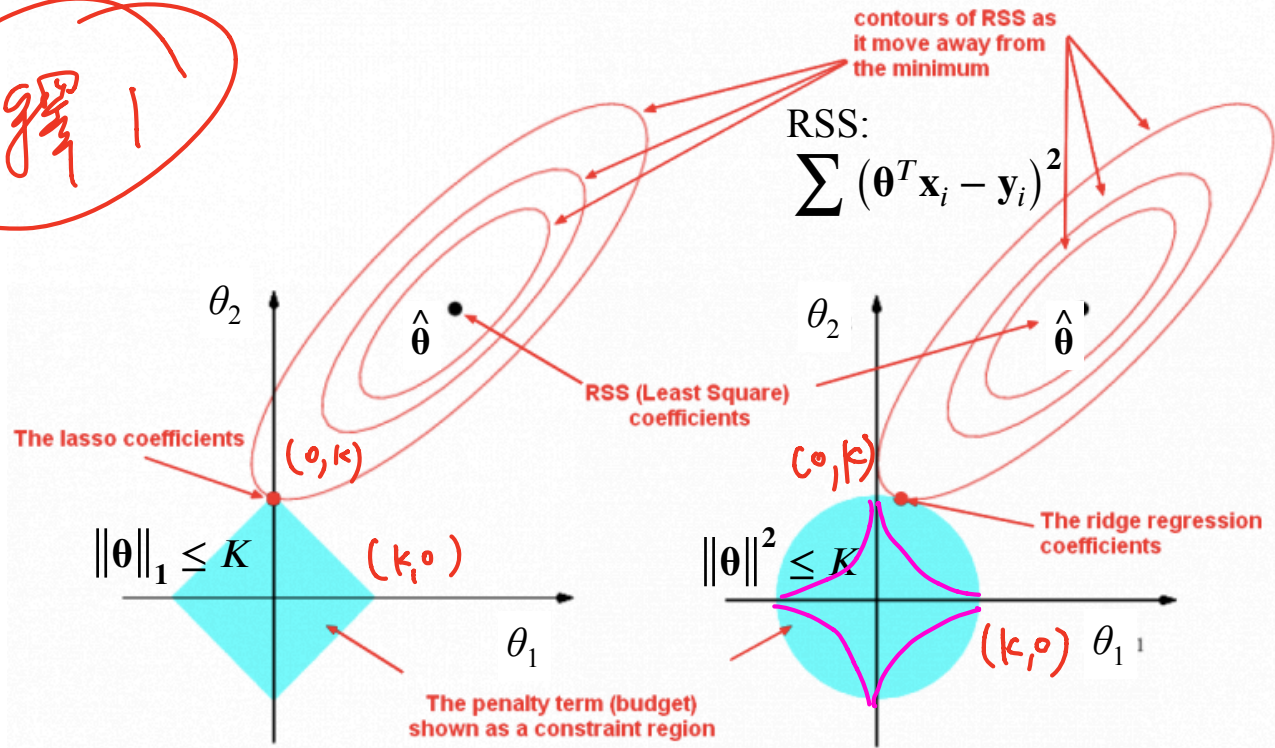
# Geometry of Ridge, Lasso, and Elastic net

Elastic net: $\lambda_1 = \lambda_2 = 0.5$

# Why Lasso zeros coefficients

解釋 I



contours of RSS as it move away from the minimum

RSS:
$$\sum \left( \theta^T \mathbf{x}_i - \mathbf{y}_i \right)^2$$

$\theta_2$

$\hat{\theta}$

RSS (Least Square) coefficients

The lasso coefficients

$(0, k)$

$\|\theta\|_1 \leq K$

$(k, 0)$

$\theta_1$

The penalty term (budget) shown as a constraint region

**LASSO**

$|\theta_1| + |\theta_2| \leq k$

$\theta_2$

$\hat{\theta}$

$(0, k)$

The ridge regression coefficients

$\|\theta\|^2 \leq K$

$(k, 0)$  $\theta_1$

**RIDGE REGRESSION**

$\theta_1^2 + \theta_2^2 \leq k$

# Why Lasso zeros coefficients? – a numerical explanation

- Ridge

  - $$\theta := \underset{\theta}{\mathrm{argmin}} \left[ \frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 + \lambda \|\theta\|^2 \right]$$

    - Changing $\theta_j$ from 2 to 1 reduces the cost by $\lambda(2^2 - 1^2) = 3\lambda$
    - Changing $\theta_j$ from 1 to 0 reduces the cost by $\lambda(1^2 - 0^2) = \lambda$

- Lasso

  - $$\theta := \underset{\theta}{\mathrm{argmin}} \left[ \frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 + \lambda \|\theta\|_1 \right]$$

    - Changing $\theta_j$ from from 2 to 1 reduces the cost by $\lambda(2 - 1) = \lambda$
    - Changing $\theta_j$ from from 1 to 0 reduces the cost by $\lambda(1 - 0) = \lambda$

➔   Ridge tends to shrink large coefficients to smaller ones, but not to shrink small coefficients to zero

# Generalize the regularization term

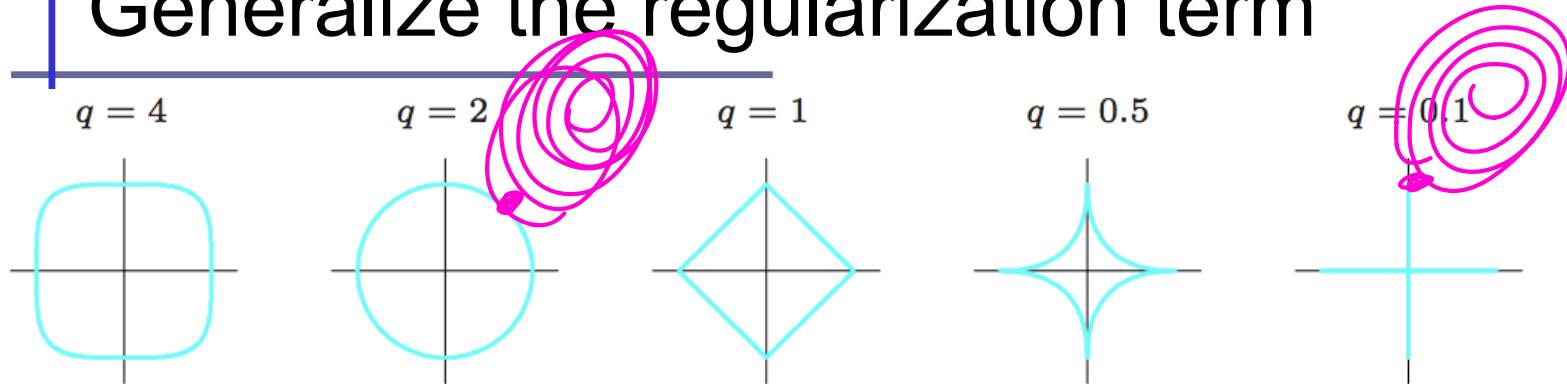| $q = 4$ | $q = 2$ | $q = 1$ | $q = 0.5$ | $q = 0.1$ |
|---------|---------|---------|-----------|-----------|

$$\boldsymbol{\theta} := \operatorname*{argmin}_{\boldsymbol{\theta}} \left[ \frac{1}{2} \sum \left( \hat{y}_i - y_i \right)^2 + \lambda \|\boldsymbol{\theta}\|^q \right]$$

$q \geq 0$

- *q*=1: Lasso
- *q*=2: Ridge regression
- A smaller *q* tends to shrink the coefficients

# Gradient descent for ridge regression

- Gradient descent (GD)
  - $\boldsymbol{\theta}^{(k+1)} := \boldsymbol{\theta}^{(k)} - \alpha \nabla J\left(\boldsymbol{\theta}^{(k)}\right)^T$

- Ridge regression
  - $J(\boldsymbol{\theta}) = \dfrac{1}{2}\left(\hat{\mathbf{y}} - \mathbf{y}\right)^T\left(\hat{\mathbf{y}} - \mathbf{y}\right) + \dfrac{\lambda}{2}\|\boldsymbol{\theta}\|^2$

    $= \dfrac{1}{2}\left(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\right)^T\left(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\right) + \dfrac{\lambda}{2}\boldsymbol{\theta}^T\boldsymbol{\theta}$

  - $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \left(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\right)^T\mathbf{X} + \lambda\boldsymbol{\theta}^T$
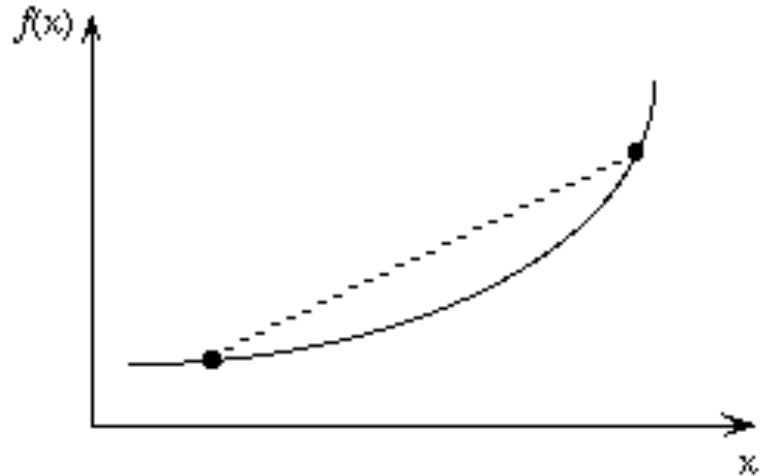
  - $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})^T = \mathbf{X}^T\left(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\right) + \lambda\boldsymbol{\theta}$

# Gradient descent for ridge regression (pseudo code)

```
1   k := 0
```

2   Initialize $\boldsymbol{\theta}^{(k)}$

3   while not converge:

4   $\quad g := \mathbf{X}^T\big(\mathbf{X}\boldsymbol{\theta}^{(k)} - \mathbf{y}\big) + \lambda\boldsymbol{\theta}^{(k)}$

5   $\quad \boldsymbol{\theta}^{(k+1)} := \boldsymbol{\theta}^{(k)} - \alpha g$
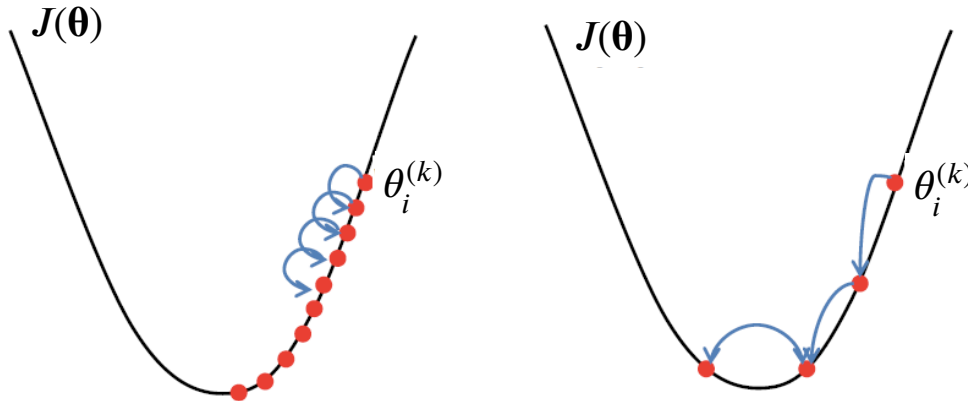
6   $\quad k := k+1$

# Convex function

- Convex function
  - Local minimum is global minimum
- The least-square linear regression objective and the Ridge regression objective are both convex functions
  - Gradient descent will find the unique minimum (or very close, depending on the step size $\alpha$

# Effect of the step size $\alpha$

- Small $\alpha$: Long computation time

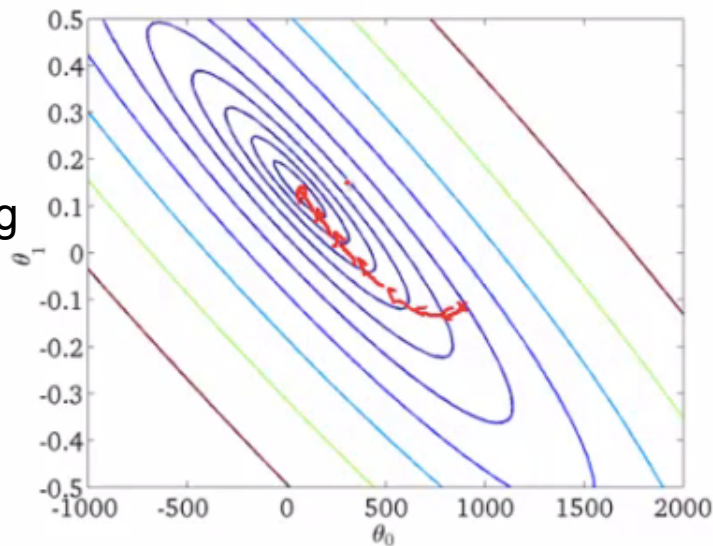- Large $\alpha$: may not reach the minimum



- A common practice
  - Gradually shrink the value of the step size

# Stochastic gradient descent (SGD) – motivation

- In GD, the gradient is computed by

  – $g := \mathbf{X}^T(\mathbf{X}\boldsymbol{\theta}^{(k)} - \mathbf{y}) + \lambda\boldsymbol{\theta}^{(k)}$

  – ***The entire training set*** is examined at each step

  – ***Very slow*** when the # of training instances is large!



```
1   k := 0
2   Initialize θ(k)
3   while not converge:
4       g := X^T(Xθ^(k) − y) + λθ^(k)
5       θ^(k+1) := θ^(k) − αg
6       k := k+1
```

# Stochastic gradient descent

- Optimize *__one example__* at a time

```
1   k := 0
2   Initialize θ(k)
3   while not converge:
4       g := x_i(x_i^T θ(k) − y_i) + λθ(k)
5       θ(k+1) := θ(k) − αg
6       k := k+1
7       Get next data instance i
```
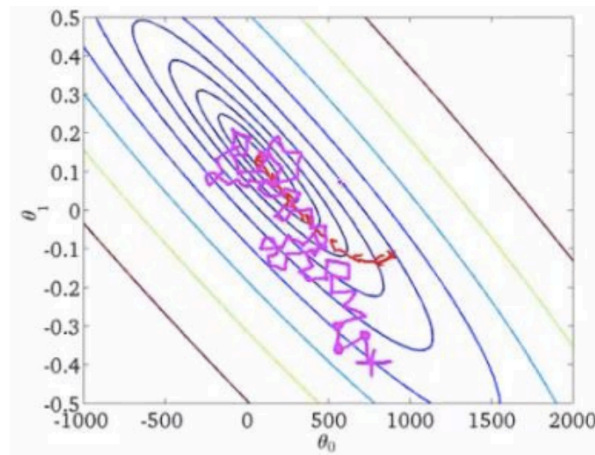
**Only one training instance is examined at each step**

- Why does SGD work?

$$- E\left(\nabla J_i(\boldsymbol{\theta})\right) = \frac{1}{n}\sum \nabla J_i(\boldsymbol{\theta}) = \nabla J(\boldsymbol{\theta})$$

# Gradient descent (GD) vs stochastic gradient descent (SGD)

- # steps
  - GD: fewer steps
  - SGD: more steps
- Computation of each step
  - GD: look through **_all_** the training instances
  - SGD: look only **_one_** training instance

# Pros and cons of SGD

- Pros
  - When the training data is large with some (near) redundant instances, SGD is usually much faster to converge than GD
  - Supports online learning
- Cons
  - Tends to bouncing around the minimum

# Mini-batch gradient descent

- Optimize *__few examples__* at a time

1   $k$ := 0

2   Initialize $\boldsymbol{\theta}^{(k)}$

3   while not converge:

4    $g := \mathbf{x}_{i:j}^{T}\left( \mathbf{x}_{i:j}\boldsymbol{\theta}^{(k)} - \mathbf{y}_{i:j} \right) + \lambda\boldsymbol{\theta}^{(k)}$

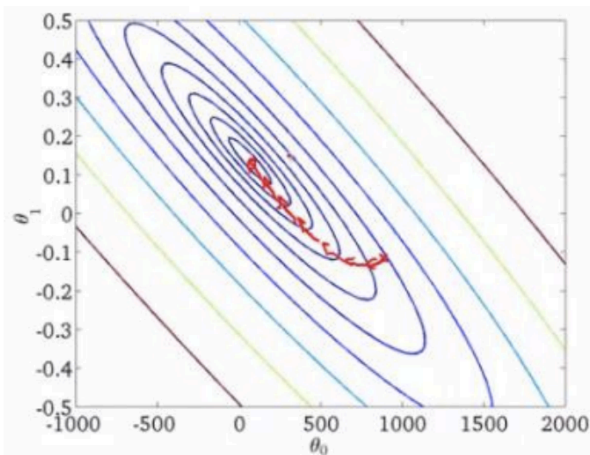5    $\boldsymbol{\theta}^{(k+1)} := \boldsymbol{\theta}^{(k)} - \alpha g$

6    $k$ := $k$+1

7    Get next batch $(i,i+1,…,j)$

Training instances $(i, i+1, …, j)$ are examined

# Gradient descent/stochastic gradient descent/mini-batch gradient descent

- All of them iteratively update the parameters such that the target function gradually becomes smaller
- If we have $n$ training instances
  - (Batch) gradient descent: every parameter update requires seeing **_all_** training instances once
  - Stochastic gradient descent: every parameter update requires seeing **_one_** of $n$ training instances
  - Mini-batch: every parameter update requires seeing **_b_** training instances (if batch size = $b$)

# Locally weighted linear regression

- Linear regression
  - Method
    - Assume $\hat{y} = \theta_0 + \theta_1 x$
    - Find $\theta_0$ and $\theta_1$ to minimize $\sum_i \left( y_i - \hat{y}_i \right)^2$
    - Given a new instance $x^*$, we predict the corresponding $y^*$ as $\theta_0 + \theta_1 x^*$
  - Once training finished, we get $\theta_0$ and $\theta_1$, and training data can be removed
  - However, can only model linear relationship
- Locally weighted linear regression
  - Model the "nonlinear" data by weighting errors differently

# Locally weighted linear regression

- Method
  - Given a new instance $x^*$, we define the weight of each training error by
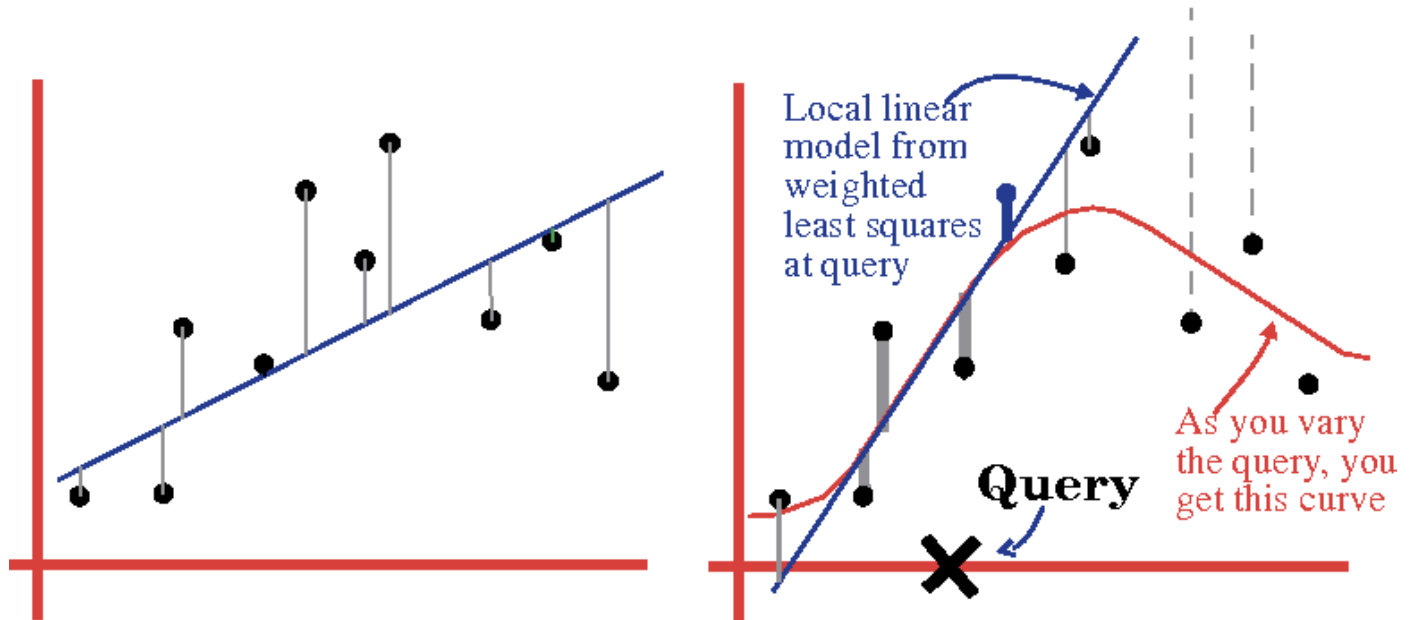    - $$w_i = \exp\left(\frac{(x_i - x^*)^2}{2\sigma^2}\right)$$
      - $\sigma$ is a hyperparameter; we set to 1 for simplicity
      - If $x_i \approx x^*$, $w_i \approx 1$
      - If $x_i$ and $x^*$ is far, $w_i \approx 0$
  - Find $\theta_0$ and $\theta_1$ to minimize
    - $$\sum_i w_i \left(y_i - \hat{y}_i\right)^2, \text{ where } \hat{y}_i = \theta_0 + \theta_1 x_i$$
- Adv: can fit nonlinear curve
- Disadv: need to "train" the model for every prediction

# LR vs locally weighted LR



Local linear model from weighted least squares at query

As you vary the query, you get this curve

**Query**

# Newton's method for optimization (1/2)

- Let $\ell(\theta) = \sum_i \left( y_i - \hat{y}_i \right)^2$, we want to

  find $\theta$ to minimize $\ell(\theta)$
  - The same as find $\theta$ such that $\ell'(\theta) = 0$
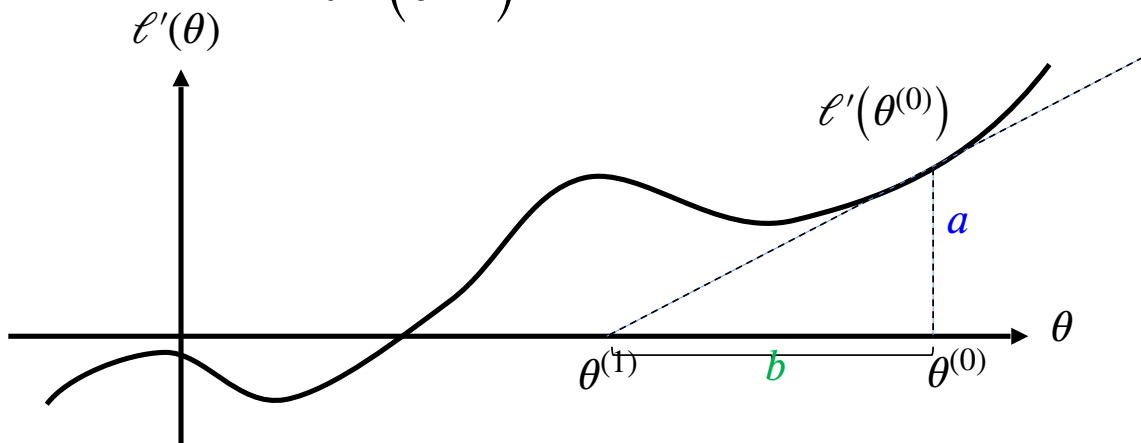
# Newton's method for optimization (2/2)

- Initial guess: $\theta^{(0)}$

- The slope of the tangent line at $(\theta^{(0)}, \ \ell'(\theta^{(0)})$ is $\ell''(\theta^{(0)})$

  - I.e., $\ell''(\theta^{(0)}) = \dfrac{a}{b} = \dfrac{\ell'(\theta^{(0)})}{\theta^{(0)} - \theta^{(1)}}$

- So, $\theta^{(1)} = \theta^{(0)} - \dfrac{\ell'(\theta^{(0)})}{\ell''(\theta^{(0)})}$

# Newton's mthod vs SGD

- Newton's method usually requires fewer steps
- However, the cost of each step is usually large
  - If $\theta$ is a scalar, using Newton's method requires to compute $f''(\theta)$
  - If $\boldsymbol{\theta}$ is a vector, using Newton's method requires to compute the Hessian matrix

# Evaluation (regression)

# Evaluating regression models

- **Mean Square Error (MSE)**

  $$MSE(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\sum \left( y_i - \hat{y}_i \right)^2}{n}$$

  – Similar metric: Root Mean Squared Error (RMSE)
  – Criticism:
    - Can only be compared between models whose errors are measured in the same unit
    - Tend to be influenced by extreme values

- Mean Absolute Error (MAE)

  $$MAE(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\sum |y_i - \hat{y}_i|}{n}$$

  – Usually smaller than MSE
  – Criticism:
    - Can only be compared between models whose errors are measured in the same unit

# Evaluating regression models

- Median Absolute Error (MedAE)
  - $MedAE(\mathbf{y}, \hat{\mathbf{y}}) = \text{median}(|y_1 - \hat{y}_1|, |y_2 - \hat{y}_2|, \ldots, |y_n - \hat{y}_n|)$
  - Robust to the extreme values
  - Criticism:
    - Can only be compared between models whose errors are measured in the same unit

- *$R^2$ score (coefficient of determination)*

  - $$R^2\left(\mathbf{y}, \hat{\mathbf{y}}\right) = 1 - \frac{\sum \left(y_i - \hat{y}_i\right)^2}{\sum \left(y_i - \bar{y}\right)^2}$$

  - The scale of the model is more intuitive
    - Best result: 1
    - Can be arbitrarily worse (very negative)

# A closer look of $R^2$ score (coefficient of determination)

- Mean model
  - A naïve regression function – a constant function that always outputs $\bar{y}$ the mean of the target variable in the training dataset

- $$R^2\left(\mathbf{y},\, \hat{\mathbf{y}}\right) = 1 - \frac{\frac{1}{n}\sum\left(y_i - \hat{y}_i\right)^2}{\frac{1}{n}\sum\left(y_i - \bar{y}\right)^2} = 1 - \frac{\sum\left(y_i - \hat{y}_i\right)^2}{\sum\left(y_i - \bar{y}\right)^2}$$

  - $\sum\left(y_i - \hat{y}_i\right)^2$ : sum of square error (SSE)

  - $\sum\left(y_i - \bar{y}\right)^2$ : sum of square total (SST)
    - i.e., SSE of the mean model

- Max value: 1
  - When SSE=0 (all predictions are correct)
- Min value: $-\infty$
  - SSE can be arbitrarily worse

- $$R^2\left(\mathbf{y},\, \hat{\mathbf{y}}\right) = 0$$

  - The performance of the prediction is the same as the mean model

## $R$ (correlation coefficient) vs $R^2$ (coefficient of determination)

- $\left[R\left(y, \hat{y}\right)\right]^2$ equals $R^2\left(y, \hat{y}\right)$ when

  – Applying the linear regression model, i.e., $\hat{y} = X\theta$, where $\theta = \left(X^T X\right)^{-1} X^T y$, and

  – Evaluating $R\left(y, \hat{y}\right)$ and $R^2\left(y, \hat{y}\right)$ on the training data

- In this case, $0 \leq R^2\left(y, \hat{y}\right) \leq 1$

# Conclusion (1/3)

- Linear regression
  - $\hat{y}_i = \theta_0 + \theta_1 x_{i,1} + \theta_2 x_{i,2} + \ldots + \theta_d x_{i,d}$
  - Find good $\theta_0$, $\theta_1$, ..., $\theta_d$ such that $\hat{y}_i \approx y_i \; \forall i$
    - Commonly used distance between $\hat{y}_i$ and $y_i$: RSS

- Linear regression with constraints on **θ**
  - Prevent overfitting
  - Limiting **θ** by L2-norm: Ridge regression
  - Limiting **θ** by L1-norm: Lasso
    - Feature selection
  - Limiting **θ** by both L1-norm and L2-norm: Elastic net

# Conclusion (2/3)

- Convex function
  - Local minimum equals global minimum
  - Iteratively adjust $\boldsymbol{\theta}$ to reduce $J(\boldsymbol{\theta})$
    - $\boldsymbol{\theta}^{(k+1)} := \boldsymbol{\theta}^{(k)} - \alpha \boldsymbol{g}$
- For ridge regression
  - Gradient descent
    - $\boldsymbol{g} := \mathbf{X}^T\big(\mathbf{X}\boldsymbol{\theta}^{(k)} - \mathbf{y}\big) + \lambda\boldsymbol{\theta}^{(k)}$
  - Stochastic gradient descent
    - $\boldsymbol{g} := \mathbf{x}_i\big(\mathbf{x}_i^T\boldsymbol{\theta}^{(k)} - \mathbf{y}_i\big) + \lambda\boldsymbol{\theta}^{(k)}$
  - Mini-batch gradient descent
    - $\boldsymbol{g} := \mathbf{x}_{i:j}^T\big(\mathbf{x}_{i:j}\boldsymbol{\theta}^{(k)} - \mathbf{y}_{i:j}\big) + \lambda\boldsymbol{\theta}^{(k)}$

# Conclusion (3/3)

- Common metrics to evaluate regression models
  - Mean square error
  - $R^2$ score

# Quiz

- If the training data contain millions of features, should we use gradient descent or closed form solution for training?

- Can gradient descent get stuck in a local minimum (which is not a global minimum) when training a linear regression model?

- Why would one use Ridge regression instead of plain linear regression?

- If the $R^2$ score on test data is negative, what does this mean?

- If the $R^2$ score on training data is negative, what does this mean?