

STA141C: Big Data & High Performance Statistical Computing

Lecture 7: Linear Regression, Linear System Solvers

Cho-Jui Hsieh
UC Davis

May 9, 2017

Linear Regression

Regression

- Input: training data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ and corresponding outputs $y_1, y_2, \dots, y_n \in \mathbb{R}$
- Training: compute a function f such that $f(\mathbf{x}_i) \approx y_i$ for all i
- Prediction: given a testing sample $\tilde{\mathbf{x}}$, predict the output as $f(\tilde{\mathbf{x}})$
- Examples:
 - Income, number of children \Rightarrow Consumer spending
 - Processes, memory \Rightarrow Power consumption
 - Financial reports \Rightarrow Risk
 - Atmospheric conditions \Rightarrow Precipitation

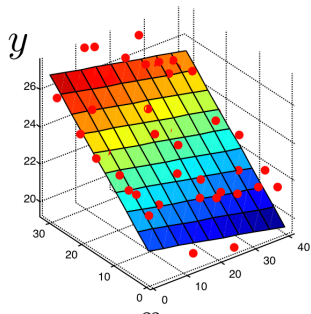
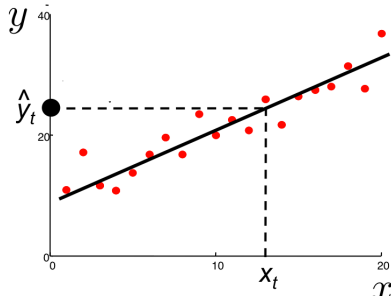
Linear Regression

- Assume $f(\cdot)$ is a linear function parameterized by $\mathbf{w} \in \mathbb{R}^d$:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

- Training: compute the model $\mathbf{w} \in \mathbb{R}^d$ such that $\mathbf{w}^T \mathbf{x}_i \approx y_i$ for all i
- Prediction: given a testing sample $\tilde{\mathbf{x}}$, the prediction value is $\mathbf{w}^T \tilde{\mathbf{x}}$
- How to find \mathbf{w} ?

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$



Linear Regression: probability interpretation

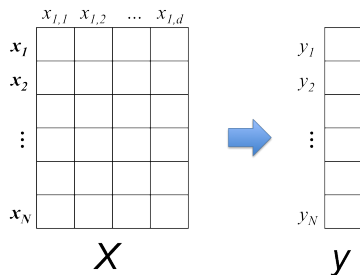
- Assume the data is generated from the probability model:

$$y_i \sim \mathbf{w}^T \mathbf{x}_i + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, 1)$$

- Maximum likelihood estimator:

$$\begin{aligned} \mathbf{w}^* &= \underset{\mathbf{w}}{\operatorname{argmax}} \log P(y_1, \dots, y_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{w}) \\ &= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^n \log P(y_i \mid \mathbf{x}_i, \mathbf{w}) \\ &= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi}} e^{-(\mathbf{w}^T \mathbf{x}_i - y_i)^2 / 2} \right) \\ &= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^n -\frac{1}{2} (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \text{constant} \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 \end{aligned}$$

Linear Regression: written as a matrix form



- Linear regression: $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$
- Matrix form: let $X \in \mathbb{R}^{n \times d}$ be the matrix where the i -th row is \mathbf{x}_i , $\mathbf{y} = [y_1, \dots, y_n]^T$, then linear regression can be written as

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

Solving Linear Regression

- Minimize the sum of squared error $J(\mathbf{w})$

$$\begin{aligned}J(\mathbf{w}) &= \frac{1}{2} \|X\mathbf{w} - \mathbf{y}\|^2 \\ &= \frac{1}{2} (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) \\ &= \frac{1}{2} \mathbf{w}^T X^T X \mathbf{w} - \mathbf{y}^T X \mathbf{w} + \frac{1}{2} \mathbf{y}^T \mathbf{y}\end{aligned}$$

- Derivative: $\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}) = X^T X \mathbf{w} - X^T \mathbf{y}$
- Setting the derivative equal to zero gives the **normal equation**

$$X^T X \mathbf{w}^* = X^T \mathbf{y}$$

- Therefore, $\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}$

Solving Linear Regression

- Minimize the sum of squared error $J(\mathbf{w})$

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \|X\mathbf{w} - \mathbf{y}\|^2 \\ &= \frac{1}{2} (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) \\ &= \frac{1}{2} \mathbf{w}^T X^T X \mathbf{w} - \mathbf{y}^T X \mathbf{w} + \frac{1}{2} \mathbf{y}^T \mathbf{y} \end{aligned}$$

- Derivative: $\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}) = X^T X \mathbf{w} - X^T \mathbf{y}$
- Setting the derivative equal to zero gives the **normal equation**

$$X^T X \mathbf{w}^* = X^T \mathbf{y}$$

- Therefore, $\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}$
but $X^T X$ may be non-invertible ...

Linear System Solver

- Linear System: given $A \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, solve (find a $\mathbf{x} \in \mathbb{R}^n$ such that)

$$A\mathbf{x} = \mathbf{b}$$

- Three cases:
 - A is invertible ($m = n$ and full rank)
 \Rightarrow unique solution $\mathbf{x} = A^{-1}\mathbf{b}$

Linear System Solver

- Linear System: given $A \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, solve (find a $\mathbf{x} \in \mathbb{R}^n$ such that)

$$A\mathbf{x} = \mathbf{b}$$

- Three cases:
 - A is invertible ($m = n$ and full rank)
 \Rightarrow unique solution $\mathbf{x} = A^{-1}\mathbf{b}$
 - **Under-determined system**: $\text{rank}(A) = m$ but $n > m$
 \Rightarrow multiple solutions, usually want to find the “least norm” solution.

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_2^2 \quad \text{s.t.} \quad A\mathbf{x} = \mathbf{b}.$$

Linear System Solver

- Linear System: given $A \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, solve (find a $\mathbf{x} \in \mathbb{R}^n$ such that)

$$A\mathbf{x} = \mathbf{b}$$

- Three cases:
 - A is invertible ($m = n$ and full rank)
 - \Rightarrow unique solution $\mathbf{x} = A^{-1}\mathbf{b}$
 - **Under-determined system**: $\text{rank}(A) = m$ but $n > m$
 - \Rightarrow multiple solutions, usually want to find the “least norm” solution.

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_2^2 \quad \text{s.t.} \quad A\mathbf{x} = \mathbf{b}.$$

- **Over-determined system**: $\text{rank}(A) < m$
 - \Rightarrow (usually) no solution
 - \Rightarrow Output $\mathbf{x}^* = \arg \min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|^2$

Linear System Solver

- Linear System: given $A \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, solve (find a $\mathbf{x} \in \mathbb{R}^n$ such that)

$$A\mathbf{x} = \mathbf{b}$$

- Three cases:
 - A is invertible ($m = n$ and full rank)
 \Rightarrow unique solution $\mathbf{x} = A^{-1}\mathbf{b}$
 - **Under-determined system**: $\text{rank}(A) = m$ but $n > m$
 \Rightarrow multiple solutions, usually want to find the “least norm” solution.

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_2^2 \quad \text{s.t.} \quad A\mathbf{x} = \mathbf{b}.$$

- **Over-determined system**: $\text{rank}(A) < m$
 \Rightarrow (usually) no solution
 \Rightarrow Output $\mathbf{x}^* = \arg \min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|^2$
- **Do not compute the inverse of a matrix!**
 - Numerical problem
 - Time consuming

Linear System Solver

- First case: $A \in \mathbb{R}^{m \times m}$, A is invertible (full rank)
- Use “Gaussian elimination” or equivalently “LU factorization”
- Call “`numpy.linalg.solve`”

```
>>> a = np.array([3,1], [1,2])
>>> b = np.array([9,8])
>>> x = np.linalg.solve(a, b)
>>> x
array([ 2.,  3.]
```

Linear System Solver

- Under-determined system and over-determined system: can be solved by SVD

$$A\mathbf{x} = \mathbf{b}$$

$$U\Sigma V^T \mathbf{x} = \mathbf{b}$$

$$\mathbf{x} = V\Sigma^\dagger U^T \mathbf{b} \text{ (psuedo-inverse)}$$

$$\Sigma^\dagger = \text{diag}(\sigma_1^{-1}, \sigma_2^{-1}, \dots, \sigma_k^{-1}, 0, \dots, 0)$$

(assume A has k nonzero singular values)

- Call “`numpy.linalg.lstsq`”

Linear System Solver

```
>>> a = scipy.rand(5,10)
>>> b = scipy.rand(5,1)
>>> x = numpy.linalg.lstsq(a,b)
>>> numpy.linalg.norm(a.dot(x[0]) - b)
7.4320704251928296e-16
```

```
>>> a = scipy.rand(5,3)
>>> b = scipy.rand(5,1)
>>> x = numpy.linalg.lstsq(a,b)
>>> numpy.linalg.norm(a.dot(x[0]) - b)
0.35374284817556079
```

Solve multiple linear systems

- Many times we need to solve

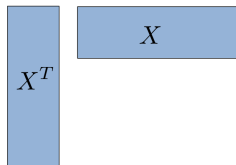
$$A\mathbf{x}_i = \mathbf{b}_i \text{ for all } i = 1, \dots, N$$

- They can be solved altogether (so only one SVD or other decomposition is needed)

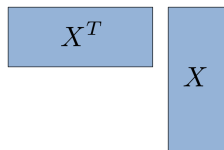
```
>>> a = scipy.rand(5,3)
>>> b = scipy.rand(5,4)
>>> x = numpy.linalg.lstsq(a,b)
>>> x[0]      ## solution for 5 linear systems
array([[ 0.15914526,  0.44365737,  0.31351924,  0.3476335 ],
       [ 0.30223114,  0.54325633,  0.22719821,  1.05852352],
       [ 0.07991735,  0.09856708,  0.08663738, -0.29111466]])
```


Solving Linear Regression

- Normal equation: $X^T X \mathbf{w}^* = X^T \mathbf{y}$
- If $X^T X$ is invertible (typically when # samples $>$ # features):
$$\mathbf{w}^* = (X^T X)^{-1} \mathbf{y}$$
- If $X^T X$ is low-rank (typically when # features $>$ # samples):
infinite number of solutions
- In general, just use “numpy.linalg.lstsq”



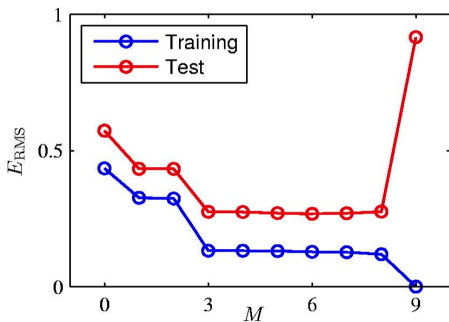
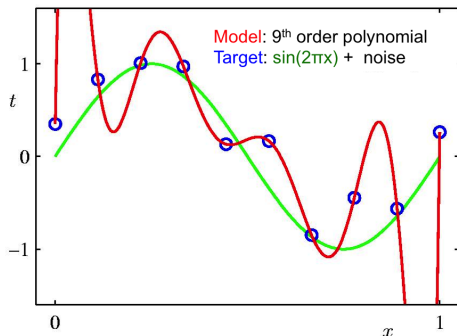
Underdetermined system
Infinite number of solutions



Overdetermined system
a unique solution

Regularized Linear Regression

Overfitting



- **Overfitting:** the model has low training error but high prediction error.
- Using too many features can lead to overfitting

Regularization to Avoid Overfitting

- Enforce the solution to have low L2-norm:

$$\operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i - y_i\|^2 \text{ s.t. } \|\mathbf{w}\|^2 \leq K$$

- Equivalent to the following problem with some λ

$$\operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i - y_i\|^2 + \lambda \|\mathbf{w}\|^2$$

Regularized Linear Regression

- Regularized Linear Regression:

$$\operatorname{argmin}_{\mathbf{w}} \|X\mathbf{w} - \mathbf{y}\|^2 + R(\mathbf{w})$$

$R(\mathbf{w})$: regularization

- Ridge Regression (ℓ_2 regularization):

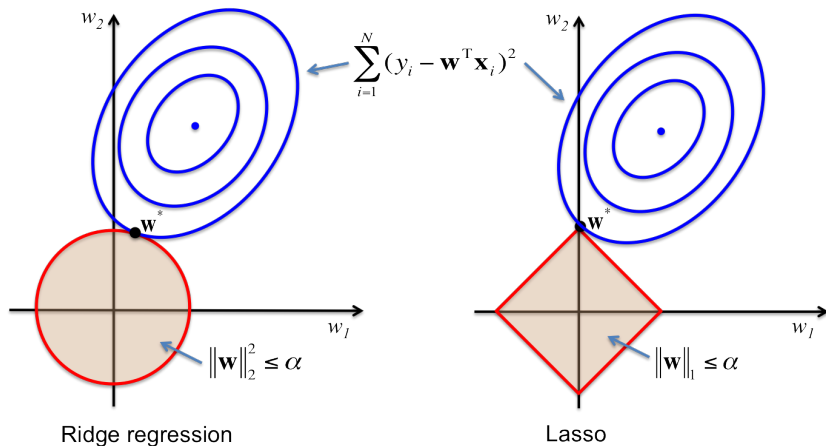
$$\operatorname{argmin}_{\mathbf{w}} \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2$$

- Lasso (ℓ_1 regularization):

$$\operatorname{argmin}_{\mathbf{w}} \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1$$

Note that $\|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|$

Regularization



Lasso: the solution is **sparse**, but no closed form solution

Ridge Regression

- Ridge regression: $\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \underbrace{\frac{1}{2} \|X\mathbf{w} - \mathbf{y}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2}_{J(\mathbf{w})}$
- Closed form solution: optimal solution \mathbf{w}^* satisfies $\nabla J(\mathbf{w}^*) = 0$:

$$\begin{aligned} X^T X \mathbf{w}^* - X^T \mathbf{y} + \lambda \mathbf{w}^* &= 0 \\ (X^T X + \lambda I) \mathbf{w}^* &= X^T \mathbf{y} \end{aligned}$$

- Optimal solution: $\mathbf{w}^* = (X^T X + \lambda I)^{-1} X^T \mathbf{y}$
- Inverse always exists because $X^T X + \lambda I$ is **positive definite**
- What's the computational complexity?

Time Complexity

- When X is dense:
 - Closed form solution requires $O(nd^2 + d^3)$ if X is dense
 - Efficient if d is very small
 - Runs forever when $d > 100,000$
- Typical case for big data applications:
 - $X \in \mathbb{R}^{n \times d}$ is sparse with large n and large d
 - How can we solve the problem?

Coming up

- Optimization

Questions?