



Machine Learning

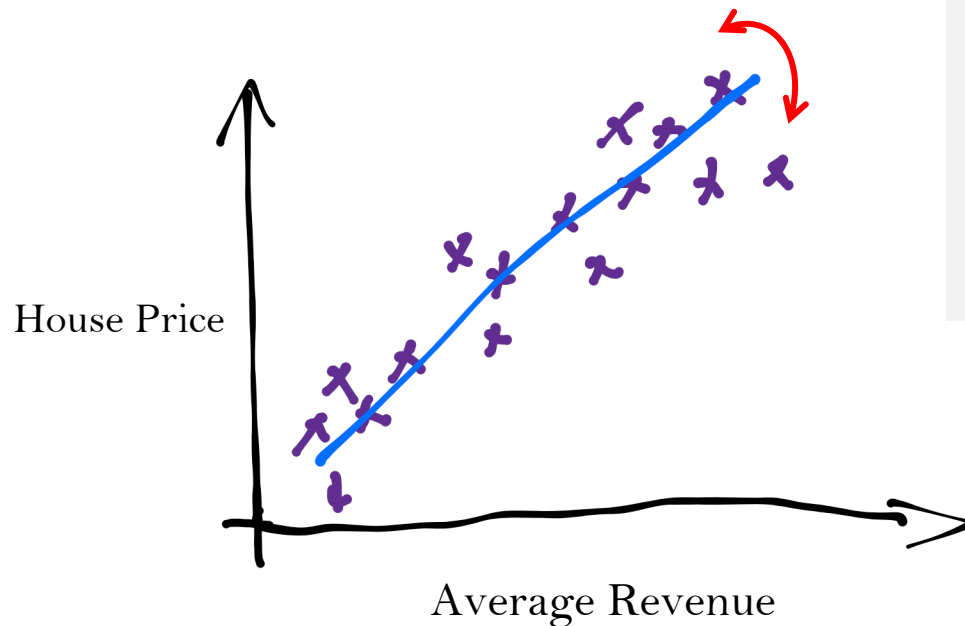
Lecture 2: Linear Regression

Fall 2023

Instructor: Xiaodong Gu



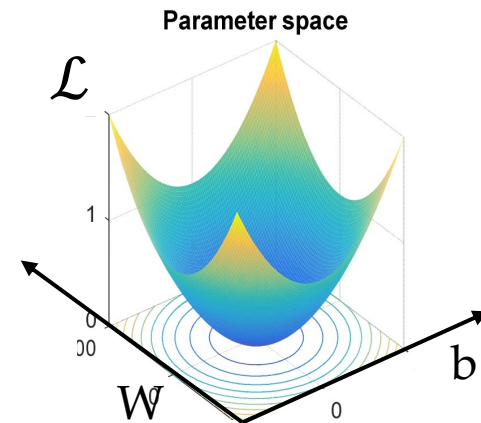
Recall: Key Elements of Machine Learning



Elements:

- #1 Data (Experience)
- #2 Model (Hypothesis)
- #3 Loss Function (Objective)
- #4 Optimization (Improve)

$$\theta^* = \operatorname{argmin}_{\theta} \mathcal{L}(\theta|\mathcal{D})$$



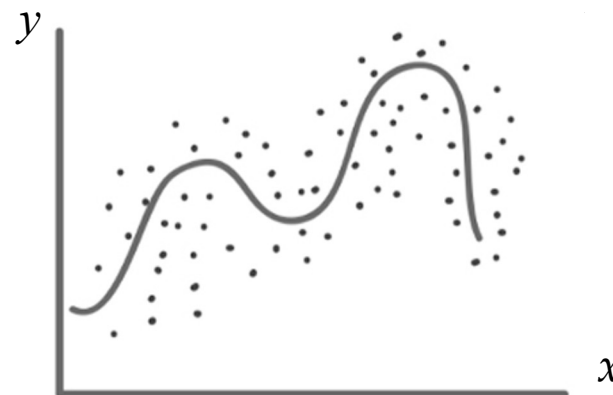


Regression in Machine Learning

Machine Learning

- **Supervised Learning**
 - Regression (✓)
 - Classification
 - ...
- Unsupervised Learning
- Reinforcement Learning

Advertisement	Sales
\$90	\$1000
\$120	\$1300
\$150	\$1800
\$100	\$1200
\$130	\$1380
\$200	??



Regression: predicts real-valued labels



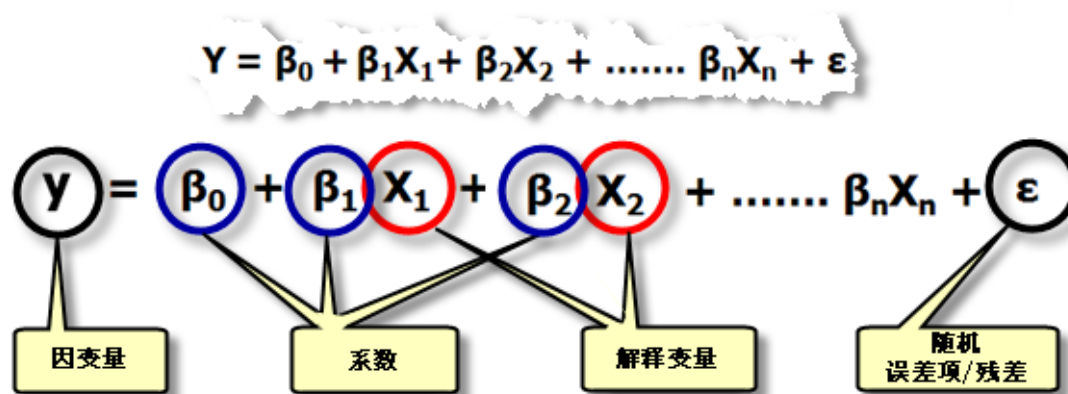
Regression

Why regression?

Regress the true value of a statistical variable through many experimentally observed values.

What is regression?

A function that describe the relationship between one **dependent** variable and a series of other (**independent**) variables.



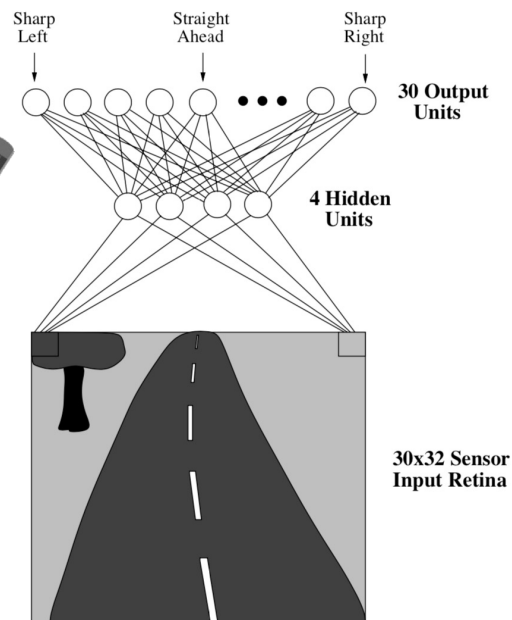
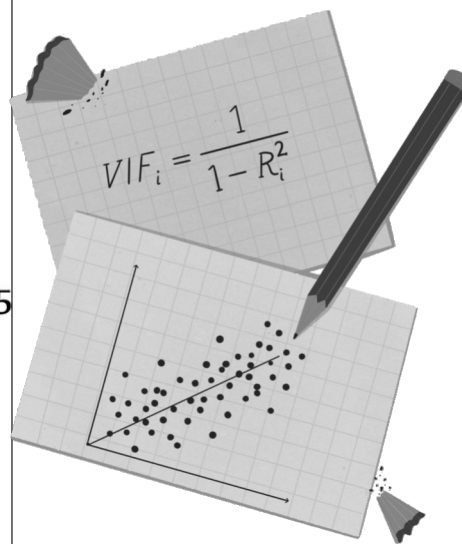
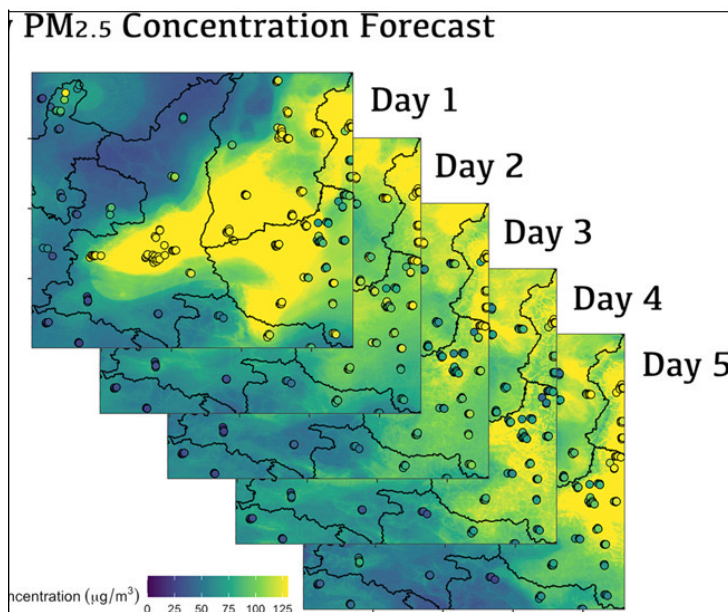
Applications of Regression

Forecasting

Factorization

Control

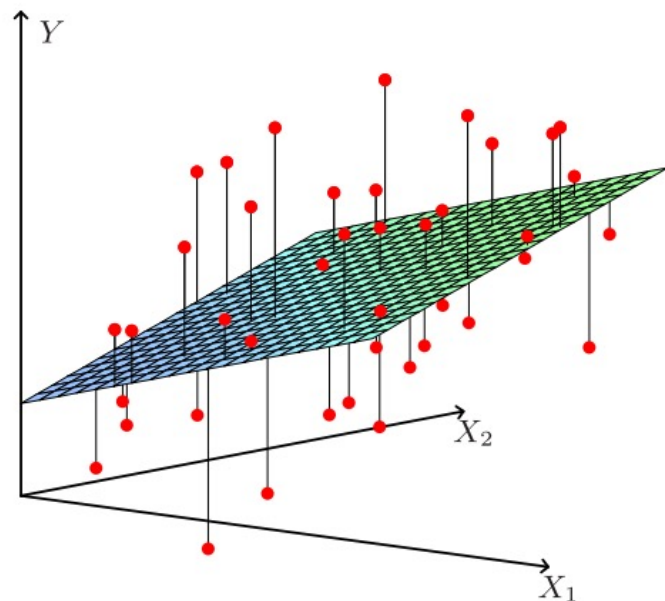
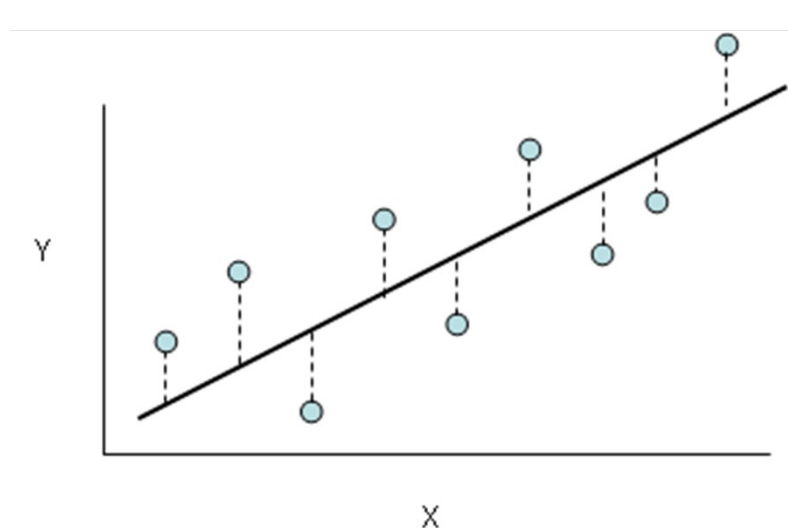
...



Linear Regression

A linear function for regression

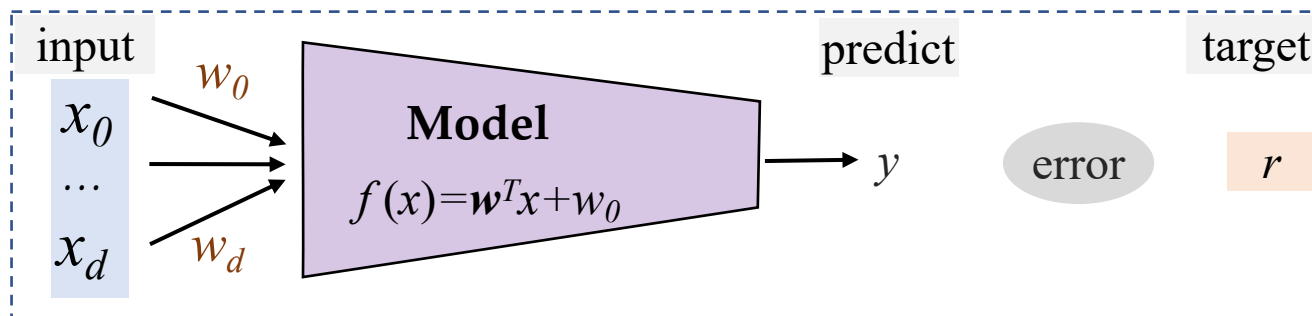
$$y = f(x) = \mathbf{w}^T \mathbf{x} + w_0$$



Linear model for regression is a $(d+1)$ -dimensional hyperplane

Model Architecture

A simple linear function.



- **Train:**
 - estimate the parameters \mathbf{w} and w_0 from data
- **Test:**
 - calculate $f(x) = \mathbf{w}^T x + w_0$.



Loss Function

- For a given input x , the model outputs a real value y . Let $r \in \mathbb{R}$ be target value, the square error is :

$$l(\mathbf{w}, w_0 | x, r) = (r - y)^2$$

- Given: $D = \{(x^{(1)}, r^{(1)}), \dots, (x^{(N)}, r^{(N)})\}$, the loss over the dataset is defined as the **mean square error (MSE)**:

$$L(\mathbf{w}, w_0 | D) = \frac{1}{2N} \sum_{\ell=1}^N (r^{(\ell)} - y^{(\ell)})^2$$

Optimization

Given: $D = \{(x^{(1)}, r^{(1)}), \dots, (x^{(N)}, r^{(N)})\}$

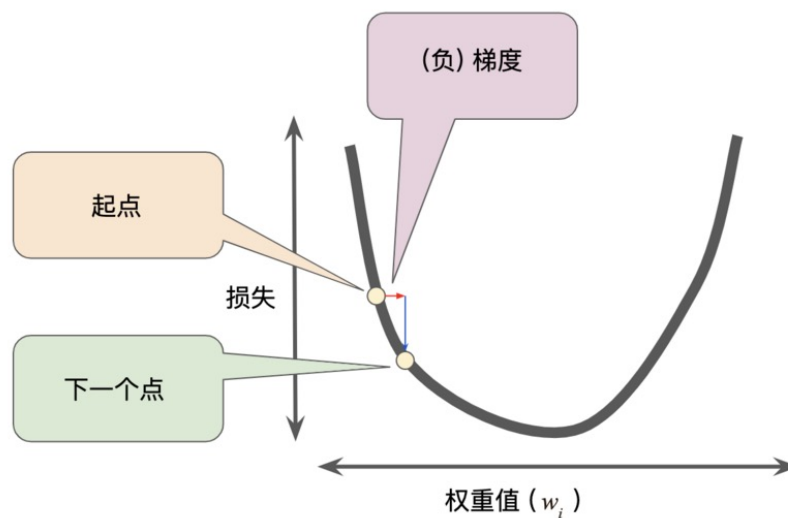
minimize the loss function using **gradient descend**:

- Goal:

$$\min_w L(\mathbf{w})$$

- Iteration:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \frac{\partial L}{\partial \mathbf{w}}$$



What is $\frac{\partial L}{\partial \mathbf{w}}$?



Optimization – Gradient Descend

$$L(\mathbf{w}, w_0 | D) = -\frac{1}{2N} \sum_{\ell=1}^N (r^{(\ell)} - y^{(\ell)})^2$$

What is $\frac{\partial L}{\partial \mathbf{w}}$?

For each w_j ($j=1, \dots, d$):

$$\frac{\partial L}{\partial w_j} = -\frac{1}{N} \sum_{\ell} \underbrace{(r^{(\ell)} - y^{(\ell)})}_{\text{Chain rule}} \frac{\partial y^{(\ell)}}{\partial w_j} = -\frac{1}{N} \sum_{\ell} (r^{(\ell)} - y^{(\ell)}) x^{(\ell)}$$

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} + \frac{1}{N} \sum_{\ell=1}^N (r^{(\ell)} - y^{(\ell)}) \mathbf{x}^{(\ell)}$$



The Algorithm

Gradient Descend for Liner Regression

Input: $D = \{(x^{(l)}, r^{(l)})\} (l=1:N)$

for $j = 0, \dots, d$

$w_j \leftarrow \text{rand}(-0.01, 0.01)$

repeat

for $j = 0, \dots, d$

$\Delta w_j \leftarrow 0$

for $l = 1, \dots, N$

$y \leftarrow 0$

for $j = 0, \dots, d$

$y \leftarrow y + w_j x_j^{(l)}$

$\Delta w_j \leftarrow \Delta w_j + (r^{(l)} - y) x_j^{(l)}$

$\Delta w_j = \Delta w_j / N$

for $j = 0, \dots, d$

$w_j \leftarrow w_j + \eta \Delta w_j$

until convergence



The Matrix Form

$$\text{Let } X = \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(N)} \end{bmatrix} = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_0^{(N)} & x_1^{(N)} & x_2^{(N)} & \dots & x_d^{(N)} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \quad \mathbf{r} = \begin{bmatrix} r^{(1)} \\ r^{(2)} \\ \vdots \\ r^{(N)} \end{bmatrix}$$

- Prediction: $\mathbf{y} = X\mathbf{w} = \begin{bmatrix} \mathbf{x}^{(1)}\mathbf{w} \\ \mathbf{x}^{(2)}\mathbf{w} \\ \vdots \\ \mathbf{x}^{(N)}\mathbf{w} \end{bmatrix}$
- Objective: $L(\mathbf{w}) = \frac{1}{2} (\mathbf{r} - \mathbf{y})^T (\mathbf{r} - \mathbf{y}) = \frac{1}{2} (\mathbf{r} - X\mathbf{w})^T (\mathbf{r} - X\mathbf{w})$



The Matrix Form

- Gradient

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = -\mathbf{X}^T(\mathbf{r} - \mathbf{X}\mathbf{w})$$

- Solution

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{X}^T(\mathbf{r} - \mathbf{X}\mathbf{w}) = 0$$

$$\Rightarrow \mathbf{X}^T \mathbf{r} = \mathbf{X}^T \mathbf{X} \mathbf{w}$$

$$\Rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{r}$$



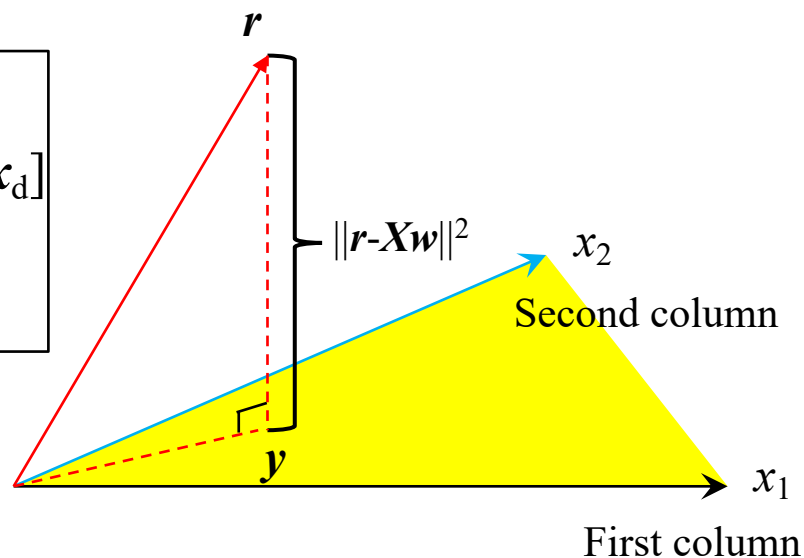
The Matrix Form

- Then the predicted values are

$$\begin{aligned} \mathbf{y} &= \underbrace{\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T}_{\mathbf{H}} \mathbf{r} \\ &= \mathbf{H}\mathbf{r} \end{aligned}$$

Geometrical Explanation

- The column vectors $[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d]$ form a subspace of \mathbb{R}^n .
- \mathbf{H} is a least square projection





Matrix Form with Regularization

Problem: $X^T X$ might be singular

When some column vectors are not independent (e.g., $\mathbf{x}_2 = 3\mathbf{x}_1$), then $X^T X$ is singular, thus $\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{r}$ cannot be directly calculated.

Solution: Regularization

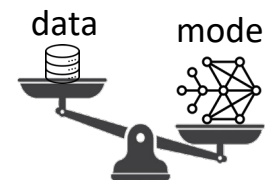
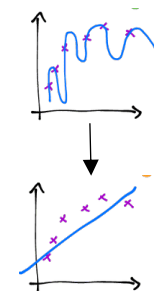
$$L(\mathbf{w}) = \frac{1}{2} (\mathbf{r} - \mathbf{y})^T (\mathbf{r} - \mathbf{y}) = \frac{1}{2} (\mathbf{r} - X\mathbf{w})^T (\mathbf{r} - X\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

New gradient: $\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = -X^T (\mathbf{r} - X\mathbf{w}) + \lambda \mathbf{w}$

New optimal solution: $\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0 \Rightarrow -X^T (\mathbf{r} - X\mathbf{w}) + \lambda \mathbf{w} = 0$

$$\Rightarrow X^T \mathbf{r} = (X^T X + \lambda \mathbf{I}) \mathbf{w}$$

$$\Rightarrow \mathbf{w}^* = (X^T X + \lambda \mathbf{I})^{-1} X^T \mathbf{r}$$



Penalty to model turns out to be data augmentation (adding data prior)

Programming Time



Tutorial:

Python

<https://colab.research.google.com/github/cs231n/cs231n.github.io/blob/master/python-colab.ipynb>

Linear regression with Python

<https://www.kaggle.com/code/sudhirnl7/linear-regression-tutorial/data?select=insurance.csv>



What's Next?



Classifications

Find a decision boundary that **maximizes the margin** between two classes.

Machine Learning

- **Supervised Learning**
 - Regression
 - Classification(✓)
 - ...
- Unsupervised Learning
- Reinforcement Learning

