

Supervised learning

DS 4400

Instructor: Ehsan Elhamifar

(A) Regression

Given input / feature / attribute vectors $\{x_i \in X\}_{i=1}^N$ and associated numeric values

$\{y_i \in \mathbb{R}\}_{i=1}^N$, find a function / hypothesis $h: X \rightarrow \mathbb{R}$ so that $h(x_i) \approx y_i, \forall i$

Ex) We want to predict the price of house

each house $x_2 = \begin{pmatrix} \# \text{bedrooms} \\ \text{square footage} \\ \text{townhouse (0) or apartment (1)} \end{pmatrix}$ has a price $y_1 = \$$

* What is our data? $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$ → training samples

* How to choose h ? Equivalently, what is the hypothesis space?

In "linear regression" $h(x_i) = \theta^T x_i$ ↑ unknown, must be learned from D. linear combination of features / attributes / variables

Eg, square footage might be more important in price than #bedrooms, etc

* How to measure closeness of $h(x_i)$ to y_i ? I.e., what objective function to choose?

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N l(h(x_i), y_i)$$

↓ loss / distance function

A popular objective / cost function is the $(l_2\text{-norm})^2$

$$l(h(x_i), y_i) = (h(x_i) - y_i)^2$$

Putting the two together, we want to minimize

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N (y_i - h_{\theta}(x_i))^2 = \sum_{i=1}^N (y_i - \theta^T x_i)^2 = J(\theta)$$

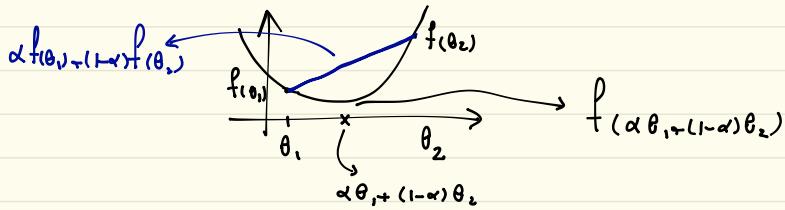
* How to solve for θ ?

Let's write this in matrix form $x_i \in \mathbb{R}^d$, $X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix} \in \mathbb{R}^{n \times d}$, $y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$

$$J(\theta) = \sum_{i=1}^n (y_i - \theta^T x_i)^2 = \|y - X\theta\|_2^2$$

* Digression $J(\theta)$ is called "convex" function iff

$$\forall \theta_1, \theta_2 \text{ and } \forall \alpha \in [0,1] \quad J(\alpha \theta_1 + (1-\alpha) \theta_2) \leq \alpha J(\theta_1) + (1-\alpha) J(\theta_2)$$



In other words, the function always lies below the segment connecting any of its two points

+ If $J(\cdot)$ is convex, then a minimizer of $J(\cdot)$ can be found by setting its derivative to zero

$$\frac{\partial J}{\partial \theta} \Big|_{\theta_{\min}} = 0 \quad (\text{gradient is zero})$$

→ Back to regression is $J(\theta) = \|Y - X\theta\|_2^2$ convex? YES

$$\left\{ \begin{array}{l} \frac{\partial \alpha^\top \theta}{\partial \theta} = \alpha \\ \frac{\partial \theta^\top A \theta}{\partial \theta} = 2A\theta \\ \|\alpha - b\|_2^2 = \|\alpha\|_2^2 - 2\alpha^\top b + \|b\|_2^2 \end{array} \right.$$

① check using first def ✓ ② $\nabla^2 J(\theta) = 2X^\top X \underset{d \times d \text{ vs } d \times 1}{\gg 0} \quad \checkmark$

How to find optimal θ (minimizer of $J(\cdot)$)?

① Analytical solution $\frac{\partial J(\theta)}{\partial \theta} \Big|_{\theta_{\min}} = 0$

$$\nabla J(\theta) = 2X^\top(X\theta - Y) = 0 \rightarrow X^\top X\theta = X^\top Y \rightarrow \theta = (X^\top X)^{-1} X^\top Y$$

↳ this should be invertible
need $N \gg d$!

Advantage: solution in one shot!

Possible drawbacks: need $N \geq d$ # points > # features

if d large need large training, also inverting $X^\top X \in \mathbb{R}^{d \times d}$ costly! $O(d^3)$

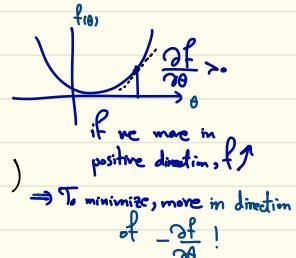
② Gradient descent algorithm

- start with $\theta^{(0)} \in \mathbb{R}^d$ (eg $\theta^{(0)} = 0$ or random vector)
- At iteration k , update θ by (until $\|\theta^{(k)} - \theta^{(k-1)}\| \leq \varepsilon$)

$$\theta^{(k)} = \theta^{(k-1)} - \gamma \frac{\partial J(\theta)}{\partial \theta} \Big|_{\theta^{(k-1)}}$$

↳ learning rate ↳ gradient of J w.r.t. θ

$$\text{Here } \theta^{(k)} = \theta^{(k-1)} - \gamma X^\top(X\theta - Y) \rightarrow \text{Batch gradient descent}$$



Advantage: less computations (still costly if $N \uparrow \rightarrow O(Nd)$)

disadvantage: can be slow (convergence rate depends on γ and $X^\top X$)

Alternatively

- start with $\theta^{(0)} \in \mathbb{R}^d$ (eg $\theta^{(0)} = 0$ or random vector)
 - At iteration k , select one point at random, update θ by
 $\xrightarrow{\text{by } x_i}$ (until $\|\theta^{(k)} - \theta^{(k-1)}\| \leq \varepsilon$)
- $$\theta^{(k)} = \theta^{(k-1)} - \gamma x_i (x_i^T \theta - y_i) \longrightarrow \text{stochastic gradient descent}$$

Advantage computationally very efficient $\rightarrow O(d)$
disadvantage can be slow

many empirical studies show it works well

Regression :

Given training samples $\{(x_1, y_1), \dots, (x_N, y_N)\}$ where $x_i \in \mathbb{R}^d$ is a feature vector ($\in X$), $y_i \in \mathbb{R}$ is the desired output (called response) ($\in Y$)

Goal find an approximate function (hypothesis) $h: X \rightarrow Y$ so that
 $h(x_i) \approx y_i, \forall i=1, N$

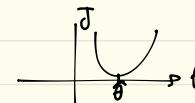
Hypothesis \Rightarrow Linear Regression $h_\theta(x) = \theta^T x$ for $\theta \in \mathbb{R}^d$, i.e., response is a linear function of unknown parameters

$$\text{Cost function } \min_{\theta} J(\theta) \triangleq \sum_{i=1}^N (y_i - h_\theta(x_i))^2 = \sum_{i=1}^N (y_i - \theta^T x_i)^2$$

$$\text{we showed } J(\theta) = \sum_{i=1}^N (y_i - \theta^T x_i) = \|Y - X\theta\|_2^2$$

$$\text{where } X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} \in \mathbb{R}^{N \times d} \text{ and } Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N$$

We saw $J(\theta)$ is convex \rightarrow allows to be minimized wrt



$$\textcircled{1} \quad \nabla J(\theta) = \frac{\partial J(\theta)}{\partial \theta} = 2 X^T (X\theta - Y) = 0 \rightarrow \hat{\theta} = (X^T X)^{-1} X^T Y$$

The strategy of
 $J(\cdot)$ is diff
wrt θ !

$$X \in \mathbb{R}^{N \times d} \Rightarrow X^T X \in \mathbb{R}^{d \times d} \text{ to be invertible needs } N \geq d$$

Adv $\hat{\theta}$ in one line (closed form)

Dis $\hat{\theta}$ computation when d is large (#pixels, words) is costly $O(d^3 + dN)$

② Gradient descent algorithm

- start with $\theta^{(0)} \in \mathbb{R}^d$ (e.g. $\theta^{(0)} = 0$ or random vector)
- At iteration K , update θ by (until $\|\theta^{(K)} - \theta^{(K-1)}\| \leq \epsilon$ or $\|\nabla J(\theta^{(K)})\| \leq \epsilon$)

$$\theta^{(k)} = \theta^{(k-1)} - \gamma \frac{\partial J(\theta)}{\partial \theta} \Big|_{\theta^{(k-1)}}$$

learning rate gradient of J w.r.t θ

Here $\theta^{(k)} = \theta^{(k-1)} - \gamma X^T (X\theta^{(k-1)} - Y)$ \rightarrow Batch gradient descent

Advantage less computations ($x\theta^{(k)}, x^T(x\theta^{(k)})$) $\rightarrow O(Nd)$ still depends on N
 disadvantage need all data at (convergence rate depends on γ and $X^T X$)
 each iteration

$N \uparrow \Rightarrow$ costly

③ Stochastic Grad Descent

- start with $\theta^{(0)} \in \mathbb{R}^d$ (e.g. $\theta^{(0)} = 0$ or random vector)
- At iteration K , select one point at random, update θ by (until $\|\theta^{(K)} - \theta^{(K-1)}\| \leq \epsilon$)

$$\theta^{(k)} = \theta^{(k-1)} - \gamma x_i (x_i^T \theta - y_i) \rightarrow \text{Replace } X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} \text{ by } X = x_i \text{ for a randomly selected point}$$

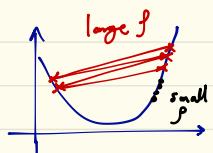
Advantage computationally very efficient $\rightarrow O(d)$

disadvantage can be slow (however, generally as slow as Batch GD)

many empirical studies show it works well

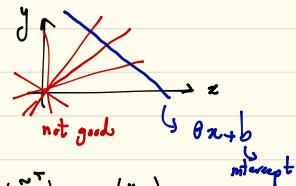
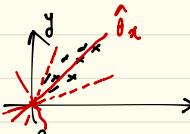
$$\sum_K f_K = \infty, \sum K^2 f_K^2 < \infty \quad \text{Robbins-Monro Condition}$$

e.g. $f_K = \frac{1}{K+0.1}$



One limitation of what we showed $h_\theta(x) = \theta^\top x$ is that it is a linear function of x , which can be restrictive

$$\text{Ex: } \left\{ (x_i \in \mathbb{R}, y_i \in \mathbb{R}) \right\}_{i=1}^n$$



To account for such affine transformations let

$$\text{feature vector be } (x_i) \rightarrow \tilde{x}_i \in \mathbb{R}^{d+1} \quad \tilde{x} = \begin{pmatrix} \tilde{x}_i^\top \\ 1 \end{pmatrix}, y = \begin{pmatrix} y_i \\ 1 \end{pmatrix}$$

$$\text{and learn } \tilde{\theta} \in \mathbb{R}^{d+1} = \begin{pmatrix} \theta \\ b \end{pmatrix}$$

Another limitation of this is its linear dependence on features/covariates $x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix} \in \mathbb{R}^d$

$$h_\theta(x_i) = \sum_{j=1}^d \theta_j x_{ij}, \quad \text{if } x_{ij} \text{ increases by } \delta_j \text{ then } h_\theta(x_i) \text{ changes by } \theta_j \delta_j!$$

This can be restrictive!

$$\text{Ex: Price Regression } x_i = \begin{pmatrix} \# bedrooms \\ \text{sq ft} \\ \text{area} \end{pmatrix}, y_i = \text{price of house},$$

if #bedrooms $1 \rightarrow 2$ I expect higher increase in price than $7 \rightarrow 8$, i.e.
"diminishing return" property \rightarrow makes sense to consider quadratic features

$$x_i = \begin{pmatrix} \# bedrooms \\ (\# bedrooms)^2 \\ \text{sq ft} \end{pmatrix}$$

$$20(2-1) - 1(2^2-1^2) = 17$$

$$20(8-7) - 1(8^2-7^2) = 5$$

1bd \rightarrow 2 bd : price increase: 17

7bd \rightarrow 8bd : price increase: 5

More generally, we can consider $\varphi(x)$ instead of x , where $\varphi: X \rightarrow H$ is a complex nonlinear function \Rightarrow Basis function expansion

$$\text{Ex: } \varphi(x) = (x, x_1^2, x_1^3, x_2, x_2^2, x_2^3, 1)^T \sim n\text{-degree polynomial}$$

Q) Is this still a linear regression?

Yes, since still the function is linear in θ $h_{\theta}(x) = \theta^T \varphi(x)$

Nonlinear in x but linear in θ \rightarrow Everything we did before, we can still do.
by using $\Phi = \begin{pmatrix} \varphi(x_1)^T \\ \varphi(x_2)^T \\ \vdots \\ \varphi(x_N)^T \end{pmatrix} \in \mathbb{R}^{N \times d'}$