

LECTURE 2: LINEAR REGRESSION AND GRADIENT DESCENT

Outline:

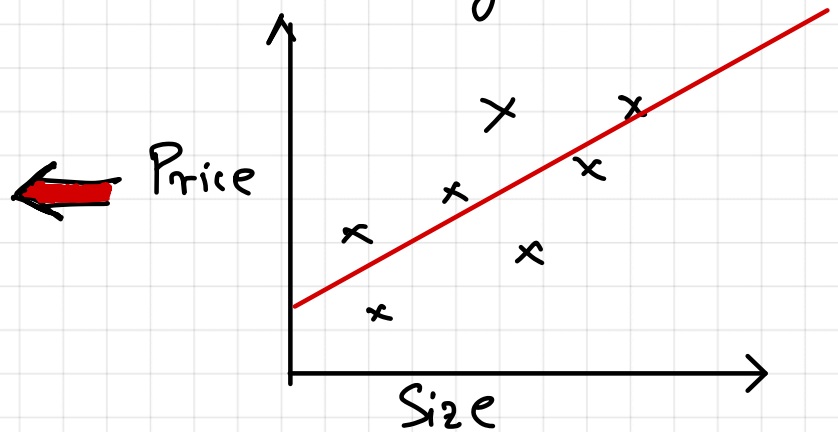
- Linear Regression
- Batch and Stochastic Gradient Descent
- Normal Equations

Supervised Learning $\begin{cases} \rightarrow \text{Regression} \\ \rightarrow \text{Classification} \end{cases}$

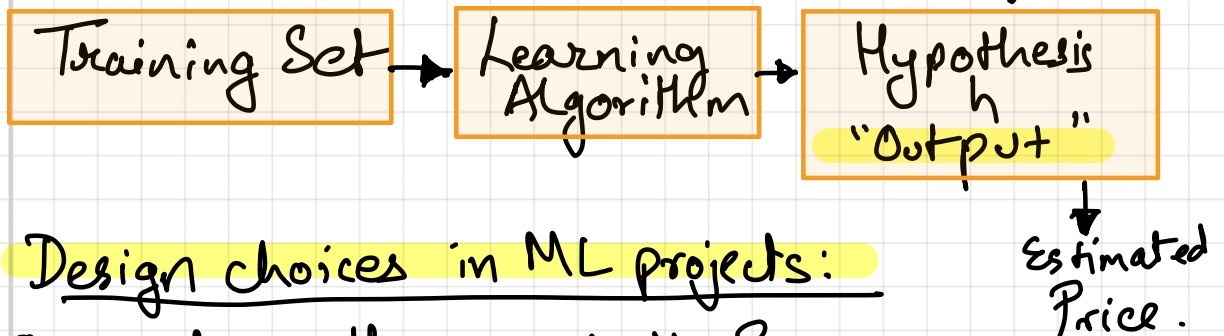
* Linear Regression:

Q. Calculating House Price using Size and

Fit a line to this data



Supervised Learning:



Design choices in ML projects:

- ① What is the workflow?
- ② What is the dataset?
- ③ How to represent the hypothesis?

Hypothesis for Linear Regression:

Let's add another feature to our problem:
of bedrooms
↓
 x_2

$$h(x) = \theta_0 + \theta_1 x \leftarrow \text{Input feature}$$

If multiple features, then: (Size in our example)

$$h(x) = \sum_{j=0}^n \theta_j x_j \quad \dots \text{ where } x_0 = 1$$

Notation:

(Used throughout the notes.)

θ : Parameters of the learning algo.

↓
Gets updated by learning algorithm

M : number of training examples/rows

X : "Inputs"/features

y : Output/target.

(x, y) : 1 training example

(x^i, y^i) : i^{th} training example

n : no. of features

→ For our problem:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$$h(x) = \sum_{j=0}^2 \theta_j x_j$$

where $x_0 = 1$

Because of x_0 & θ_0 , x and θ are $n+1$ dimensional features.

For future reference.

→ So how do you choose Theta?

Choose θ so that $h(x)$ is as close to y as possible.

So we write: $h_{\theta}(x)$ because it depends on both θ and x ,
||
 $h(x)$ although either is fine.

→ Ordinary Least Squares (OLS) Regression.

For correct/accurate/close predictions on the training set:

$$J(\theta) \Rightarrow \underset{\theta}{\text{minimize}} \frac{1}{2} \sum_{i=1}^M (h_{\theta}(x) - y)^2$$

Diagram illustrating the components of the cost function $J(\theta)$:

- Simplicity** (points to θ): for math, because we take derivative later
- Algorithm prediction** (points to $h_{\theta}(x)$)
- True value/price of the house** (points to y)

This equation is the cost function $J(\theta)$, we find parameters " θ " that minimize the cost function $J(\theta)$.

Gradient Descent

↳ Minimizes $J(\theta)$

- ① Start with some value of θ . (say $\theta = \vec{0}$)
- ② Keep changing θ to reduce $J(\theta)$.

Note:

we use

$-\alpha$
because
if we used
" $+\alpha$ " it
would move
away from
convergence.

The direction
of slope
automatically
compensates
for the "-ve"

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

← One step of Gradient Descent

↳ $j = 0, 1, 2, \dots, n$

↳ learning rate
(hyper-parameter)

$d(f(x)) \rightarrow$ defines direction of steepest descent

Partial derivative of cost function

We will assume for 1 training example,
(however it is summed over all training examples.)

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\ &= \cancel{2} \cdot \frac{1}{\cancel{2}} \cdot (h_{\theta}(x) - y) \frac{\partial (h_{\theta}(x) - y)}{\partial \theta_j} \\ &= (h_{\theta}(x) - y) \frac{\partial (\theta_0 + \theta_1 x_1 + \dots - y)}{\partial \theta_j}\end{aligned}$$

All terms will be zero, except term corresponding to $j \Rightarrow \theta_j x_j$
e.g. if $j=1$ then $\theta_1 x_1$

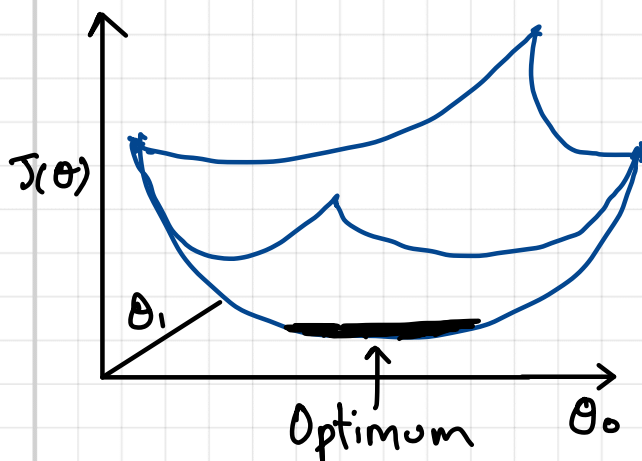
$$\therefore \frac{\partial J(\theta)}{\partial \theta_j} = (h_{\theta}(x) - y) X_j$$

One step of Gradient Descent is:

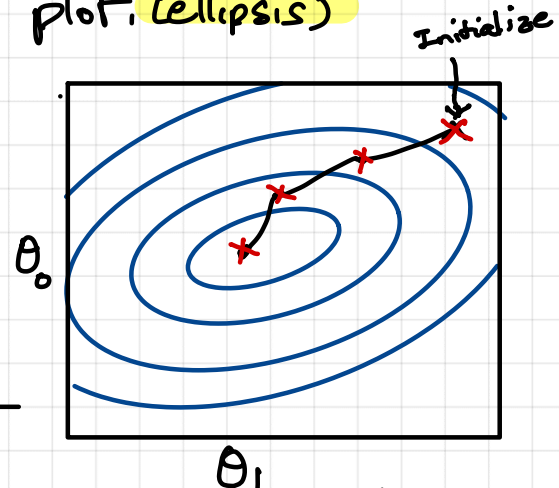
$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) X_j^{(i)}$$

For each iteration of gradient descent, carry out the update till $j=1, 2, \dots, n$,
 \Downarrow
 no. of features

→ Repeat Gradient Descent until convergence.
 $J(\theta)$ for linear regression is a quadratic function. So no local optima, only global optima. (Imagine a big bowl)



Alternatively look at the contours of the plot. (ellipses)



Step size: depends on α .

$\alpha \Rightarrow$ too large \rightarrow overshoots minima.

$\alpha \Rightarrow$ too small \rightarrow Takes too long to converge

Try a few values.

In contours, direction of steepest descent is always at 90° orthogonal to contour direction

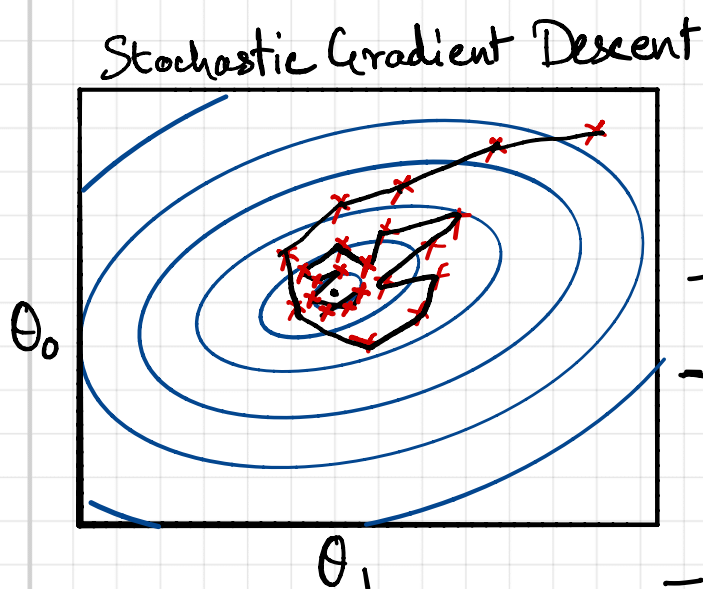
Gradient Descent

Batch Gradient Descent

- Looks at training data as one batch
- Disadvantage is that each step is very slow because we scan over entire data for each training step
- Very expensive.

Stochastic Gradient Descent

- Update parameters on every training example, instead of the full batch.
- Noisy and much more random, because your model is always oscillating between each training example.
- Never really converges
- However, much more practical in large datasets (thousands, millions)
- You end up pretty close to the global optimum, and get good enough predictions.



★ In practice with SGD, we slowly decrease learning rate over time, so you take smaller and smaller steps and end up really, really close to global optimum.

- For linear regression (only): another method to get straight to global optimum without iterative algorithm.

Normal Equations.

$J(\theta)$ → function mapping from parameters to real no.s.

Matrix
Derivative

$$\nabla_{\theta} J(\theta)$$

$$\theta \in \mathbb{R}^{n+1} \rightarrow (n+1 \text{ dimensional vector})$$

↳ Derivative of J w.r.t θ

In our problem: $\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix}$

\Downarrow

$\theta \in \mathbb{R}^{2+1}$

Matrix Derivative

e.g.

$$\nabla_A f(A) \rightarrow A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

$$f(A) = A_{11} + A_{12}^2$$

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f(A)}{\partial A_{11}} & \frac{\partial f(A)}{\partial A_{12}} \\ \frac{\partial f(A)}{\partial A_{21}} & \frac{\partial f(A)}{\partial A_{22}} \end{bmatrix} = \begin{bmatrix} 1 & 2A_{12} \\ 0 & 0 \end{bmatrix}$$

Steps: $\nabla_{\theta} J(\theta) \stackrel{\text{set}}{=} 0 \rightarrow \text{Solve for } \theta$

Trace of A = Sum of diagonal elements $= \sum A_{ii}$
($\text{tr} A$)

Properties of trace:

<p>(1) $\text{tr}(A) = \text{tr}(A^T)$</p> <p>(2) $f(A) = \text{tr} AB$ $\nabla_A f(A) = B^T$</p>	<p>(3) $\text{tr} AB = \text{tr} BA$ $\text{tr} ABC = \text{tr} CAB$ (cyclic rotation)</p> <p>(4) $\nabla_A \text{tr} AA^T C = CA + C^T A$</p>
---	---

Solve:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x^i) - y^i)^2$$

$$X = \begin{bmatrix} (x^1)^T \\ (x^2)^T \\ \vdots \\ (x^m)^T \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$= \begin{bmatrix} x^{(1)T} \theta \\ x^{(2)T} \theta \\ \vdots \\ x^{(m)T} \theta \end{bmatrix} = \begin{bmatrix} h_{\theta}(x^1) \\ h_{\theta}(x^2) \\ \vdots \\ h_{\theta}(x^m) \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

So now,

$$J(\theta) = \frac{1}{2} \underbrace{(\mathbf{X}\theta - \mathbf{y})^T}_{\downarrow} \underbrace{(\mathbf{X}\theta - \mathbf{y})}_{\downarrow} \quad \mathbf{z}\mathbf{z}^T = \sum z^2$$

$$\begin{bmatrix} h_{\theta}(x_1) - y_1 \\ \vdots \\ h_{\theta}(x_m) - y_m \end{bmatrix}$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (\mathbf{X}\theta - \mathbf{y})^T (\mathbf{X}\theta - \mathbf{y}) \\ &= \frac{1}{2} \nabla_{\theta} (\theta^T \mathbf{X}^T - \mathbf{y}^T) (\mathbf{X}\theta - \mathbf{y}) \\ &= \frac{1}{2} \nabla_{\theta} [\underbrace{\theta^T \mathbf{X}^T \mathbf{X} \theta}_{\text{}} - \underbrace{\theta^T \mathbf{X}^T \mathbf{y}}_{\text{}} - \underbrace{\mathbf{y}^T \mathbf{X} \theta}_{\text{}} - \underbrace{\mathbf{y}^T \mathbf{y}}_{\text{}}] \end{aligned}$$

Take derivative of each term

$$= \frac{1}{2} [\mathbf{X}^T \mathbf{X} \theta + \mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y}]$$

→ Refer lecture notes for derivatives

$$= \mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \mathbf{y} = \vec{0}$$

Normal equation

$$\mathbf{X}^T \mathbf{X} \theta = \mathbf{X}^T \mathbf{y}$$

$$\therefore \theta = (\mathbf{X}^T \mathbf{X}^{-1}) (\mathbf{X}^T \mathbf{y}) \rightarrow \text{One step solution for parameters}$$

→ Normal Equations.