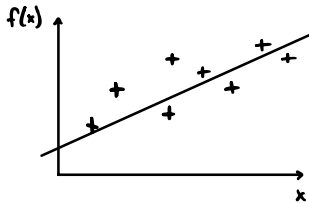


Lecture 2 - Regression

Goal: learn real valued mapping $f: \mathbb{R}^d \rightarrow \mathbb{R}$

Linear Regression



$$x = [x_1, \dots, x_d]^T, \quad f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$w = [w_1, \dots, w_d]^T, \quad f(x) = \underline{w}^T x$$

$$1\text{-dim: } y = f(x) = ax + b$$

$$2\text{-dim: } y = w_1 x_1 + w_2 x_2 + w_0$$

$$d\text{-dim: } y = w_1 x_1 + \dots + w_d x_d + w_0$$

$$= \sum_{i=1}^d w_i x_i + w_0$$

$$= w^T x + w_0$$



Homogeneous Representation

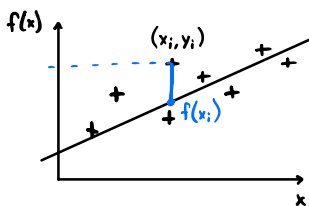
$$\tilde{w} = [w_1, \dots, w_d, w_0] = [w^T, w_0]^T$$

$$\tilde{x} = [x_1, \dots, x_d, 1] = [x^T, 1]$$

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

* this is just a way to simplify what we write by eliminating the need to separately tack on the w_0

Quantifying goodness of fit



$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$

$$\text{Residual: } r_i = y_i - w^T x_i$$

$$\text{cost: } \hat{R}(\underline{w}) = \sum_{i=1}^n r_i^2$$

* notice

↳ 1: we square each residual

↳ 2: we sum the squares

least-squares...!



Least-squares linear regression optimization

↳ given data set $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$

↳ how do we find optimal weight vector?

$$x_i \in \mathbb{R}^d$$

$$y_i \in \mathbb{R}$$

Method 1: Closed form solution

$$w^* = \arg \min_w \sum_{i=1}^n (y_i - w^T x_i)^2 \quad \text{Empirical Risk Minimizer}$$

derivation (

$$\hat{R}(w)$$

1. $R(w) = \sum_{i=1}^n (y_i - w^T x_i)^2$

2. Take gradient with respect to w and set to 0.

$$\nabla_w R(w) = -2X^T(Y - Xw)$$

$$-2X^T(Y - Xw) = 0$$

3. Rearrange...

$$X^T X w = X^T Y$$

4. Now, if $X^T X$ is invertible...

$$(X^T X)^{-1} (X^T X) w = (X^T X)^{-1} (X^T Y)$$

$$w^* = (X^T X)^{-1} X^T Y$$



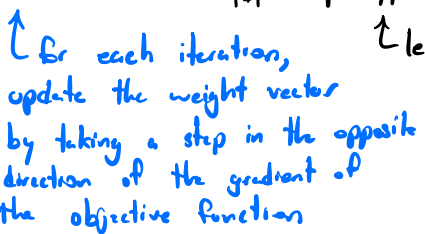
Time Complexity of $w^* = (X^T X)^{-1} X^T Y$

$$\underbrace{\begin{pmatrix} X^T X \end{pmatrix}_{n \times d \times n \times d}^{-1}}_{O(nd^2)} \uparrow \underbrace{X^T Y}_{\begin{smallmatrix} X_{n \times d} \\ Y_{n \times 1} \end{smallmatrix}}_{O(nd)} \\ \Rightarrow O(nd^2 + d^3)$$

Method 2: Optimization

The objective function $\hat{R}(w) = \sum_i (y_i - w^T x_i)^2$ is **convex!**

Gradient Descent

- Start at arbitrary $w_0 \in \mathbb{R}^d$ 
- For $t=1, 2, \dots$ do $w_{t+1} = w_t - \eta_t \nabla \hat{R}(w_t)$
 -  learning rate (e.g. $\frac{1}{2}$ for least squares)
 -  for each iteration, update the weight vector by taking a step in the opposite direction of the gradient of the objective function

Convergence

- ↳ converges to stationary point $\nabla = 0$
- ↳ optimal!
- ↳ $O(\log \frac{1}{\epsilon})$ to get within ϵ of the minimum of objective function.

Time Complexity: $O(nd)$ to compute gradient

Adaptive Step Size

- ↳ line search
- ↳ "bold driver" heuristic
 - ↳ if function \downarrow , step size \uparrow
 - ↳ if function \uparrow , step size \downarrow

Closed form vs Gradient Descent?

- computational complexity
- don't always need 100% optimal
- not all problems of closed form solution

Linear regression for polynomials

↳ we can fit **non-linear functions** via **linear regression** using non-linear features of our data

$$f(x) = \sum_i^D \omega_i \phi_i(x) \quad x \xrightarrow{\phi} \tilde{x} = \phi(x) \in \mathbb{R}^D$$

$$1\text{-d: } \phi(x) = [1, x, x^2, \dots, x^k], \quad k = D-1$$

$$2\text{-d: } \phi(x) := \phi([x_1, x_2]) = [1, x_1, x_2, x_1 x_2, x_1^2 x_2, \dots]$$

$$\vdots$$
$$p\text{-dim } \phi(x) = \mathcal{O}(p^k) \text{ vs } \mathcal{O}(k^p)$$