

14

Recap: GD.

Today: SGD Model Regularization

Gradient Descent:

$$\{x_i, y_i\} \quad i=1 \dots n$$

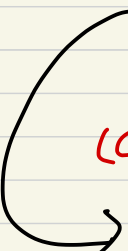
Find linear function f that $y_i \approx f(x_i)$

Alg 1: Matrix Inversion

$$y = Xw \quad w = \underbrace{(X^T X)^{-1}}_{\text{slow}} X^T y$$

Alg 2: Gradient Descent

Iterate:

$$\begin{aligned} w &\leftarrow w - \alpha \nabla L(w) \quad \text{MSE} \\ L &= \frac{1}{2} \|y - Xw\|_2^2 \\ &\quad \text{(constant)} \quad \text{assume } n=2 \end{aligned}$$

$$w_{k+1} \leftarrow w_k + \alpha_k X^T (y - Xw_k)$$

Pros: ① efficient learning time $O(ndT)$ $T = \log_{\frac{1}{\rho}} \frac{\|w^*\|_2}{\epsilon}$

$$\rho = \frac{1-\frac{L}{F}}{1+\frac{L}{F}} \quad (0 \sim 1)$$

② Simple to implement

Cons: ① possible it stuck in

Local minimum

② Need to choose α, T

↳ can be set via "line-search"

③ Requires multiple passes over data (Not Memory efficient)

④ O(n²) can be impacted by large n, d

What if dataset is too big that can't fit in Mem.

⇓
SGD

Stochastic Gradient Descent (SGD)

"Back propagation"

① How to speed up GD ?

$$\text{GD: } W_{k+1} = W_k + \alpha_k X^T (y - kW_k)$$

$$= W_k + \alpha_k \sum_{i=1}^n (y_i - \langle W_k, x_i \rangle) x_i$$

$$= W_k + \underbrace{n \alpha_k}_{\text{learning rate}} \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \langle W_k, x_i \rangle) x_i}_{\text{Average of } (y - \langle W_k, x_i \rangle) x_i}$$

Average of $(y - \langle W_k, x_i \rangle) x_i$

idea: instead of computing average over all data points, pull a subset S uniformly at random.

$$(SGD) \quad W_{k+1} = W_k + \alpha_k \sum_{i \in S} (y_i - \langle W_k, x_i \rangle) x_i$$

"Stochastic Gradient Descent"

size of $S \rightarrow$ whatever you like

(even a single point)

more points, more variance, less bias.

trade-off

SGD (single data point) 单点

(1) choose $i \in \{1, 2, \dots, n\}$
(x_i, y_i)

$$(2) \quad W_{k+1} = W_k + \alpha_k x_i (y_i - \langle W_k, x_i \rangle)$$

$\frac{1}{k}$ is choose good.

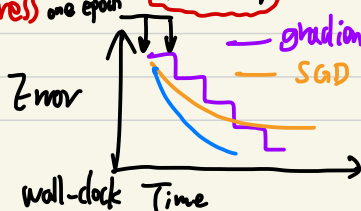
Running Time: $O(d, \#epochs)$

it won't be too small that we can't make progress

if learning rate

$$\alpha_k \propto \frac{1}{k}$$

if learning rate, then #epochs $\propto \frac{1}{\epsilon}$



— gradient descent — sth better than SGD (SVRG, SAG, ...)

SGD 其实也不是平滑的只是相比更平滑,

每个 epoch 时间更小

Model Selection

Is linear Model the right thing to do?

problems:

① Data X , label y is a nonlinear function of X

Solution: preprocess X to $[x^1, x^2, x^3, \dots, x^d, \dots, x^1]$

车比成这些, 看哪个更好

...exp(-x)

一次方, 二次方, 三次方, ... 处理

d 选得太复杂 \Rightarrow New problem: Overfitting 过拟合
正则控制 (次数据太大 --)

② Challenge in very large dataset

not all feature are relevant 不是所有 data 都相关.

eg. glucose prediction, we only need a subset of

[height, age, sugar level, ...]

③ insufficient data 不可解)

Linear Regression $W = (X^T X)^{-1} X^T Y$

$d \times d$ matrix
invertible when $\text{rank}(X^T X) = d$

$$\begin{bmatrix} | & | & & | \\ x_1 & x_2 & \dots & x_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} -x_1 - \\ -x_2 - \\ \vdots \\ -x_n - \end{bmatrix}$$

$$X^T X = \sum_{i=1}^n x_i x_i^T$$

invertible only when $n \geq d$

? 怎么知道?

#

Solution to overfitting : model selection

training dataset

$(x_1, y_1) \dots (x_n, y_n)$

pretend data unknown but follow relationship

$$y = f(x) + \epsilon_{\text{noise}}$$

"True model"

real function

Really, we care about

model we learn

$$\text{TEST MSE} = E_D (y - f(x))^2$$

distribution of data

simulated by hold-out set of training datapoints
已有所有数据的一部分

避免人为选择的影响

To avoid artifact in constructing hold-out set, repeat simulation k times

" k -fold cross validation"

见书

$$\text{TEST MSE} = E(y - f(x))^2 = E(t(x) + \epsilon - f(x))^2$$

不能用训练的数据去算 MSE,

一部分训练, 一部分 test.

$$= E(\epsilon^2) + E[(t(x) - f(x))^2]$$

$$+ 2 E(\epsilon) E(t(x) - f(x))$$

typically zero mean

真实模型 $t(x) + \epsilon$

$$= E(\epsilon^2) + E[(t(x) - f(x))^2]$$

$$E(\epsilon^2) = (E(\epsilon))^2 + \text{Var}(\epsilon)$$

$$\rightarrow E(\epsilon^2) + [E(t(x) - f(x))]^2 + \text{Var}(t(x) - f(x))$$

$$\text{TEST MSE} > \underbrace{\mathbb{E}(\epsilon^2)}_{\text{Noise}} + \underbrace{\text{Bias}^2 + \text{Variance}}_{\text{ERROR}}$$

Noise

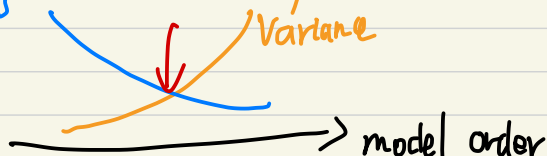
ERROR

真实模型的方差

2) Bias^2 , decrease with model complexity

3) Variance, \uparrow with model complexity

Bias²



k -fold:

将 dataset 分为 k 份, 每次循环挑一个份作为 test, 剩下

$k-1$ 份为 train, 然后挑出结果最好的作为模型

对数据做交叉处理, 得到更好结果, 因为每个参数的

大小不一样会导致权重差别很大, 实际上都一样重要.