# Lecture 6

### Gradient Descent



### Last Time:

- Machine learning, especially supervised learning
- Bias, variance, and overfitting
- Minimized an objective function, called error or cost or risk



### Today: machine learning (contd), optimization using gradient descent

- overfitting. complexity, and test sets
- gradient descent
- stochastic gradient descent

Remember Convex (bowl) like functions have 1 global minimum





### Statement of the Learning Problem

The sample must be representative of the population!

A: Empirical risk estimates in-sample risk. B: Thus the out of sample risk is also small.



 $A:R_{\mathcal{D}}(g) \; smallest \, on \, \mathcal{H}$  $B: R_{out}(g) \approx R_{\mathcal{D}}(g)$ 

### LLN: Expectations -> sample averages

$$E_p[R] = \int R(x) p(x) dx = \lim_{n o \infty} rac{1}{N} \sum_{x_i \sim x_i} e_i e_i e_i$$

### **Empirical Risk Minimization:**

$$R_{\mathcal{D}} = E_p[R] \sim rac{1}{N} \sum_{x_i \sim p} R(x_i)$$

### on training set(sample) $\mathcal{D}$ .



 $\sum_{x_i \sim p} R(x_i)$ 

### What we'd really like: population

i.e. out of sample RISK

$$R_{out}(h,y)=E_{p(x)}[R(h(x),y)]=\int dx p(x)(h(x))$$

$$\langle R_{out} 
angle = E_{p(x,y)}[R(h(x),y)] = \int dy dx \, p(x,y)$$

$$=\int dy dx p(y \mid x) p(x) R(h(x),y) = \int dx p(x) E_{p(x)} dx$$



 $(x) - y)^2 (e. g.).$ 

A(h(x),y)

 $_{(y|x)}[R(h(x),y)]$ 

- This is an average over our sampling distribution, if we had it
- What do we do?

Fit hypothesis  $h = g_{\mathcal{D}}$ , where  $\mathcal{D}$  is our training sample. Then we'd like

$$\langle R_{out} 
angle = E_{\mathcal{D}}[R_{out}(g_{\mathcal{D}},y)].$$

But:



### **Empirical Risk Minimization**

- But we only have the in-sample risk
- Furthermore its an empirical risk
- And its not even a full on empirical distribution, as N is usually quite finite

(another way of stating the LLN: the sample empirical distribution converges to the true population distribution as  $N \to \infty$ )



# What to do? TRAIN and TEST sets







**M** 207

## Is the In-Sample error small? OR FINDING DERIVATIVES



### Newton's Method



Find a zero of the first derivative.



### **Gradients and Hessians**

$$J(ar{ heta})= heta_1^2+ heta_2^2$$

Gradient: 
$$\nabla_{\bar{\theta}} (J) = \frac{\partial J}{\partial \bar{\theta}} = \begin{pmatrix} 2\theta_1 \\ 2\theta_2 \end{pmatrix}$$
  
Hessian H =  $\begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$ 

Hessian gives curvature. Why not use it?



### Gradient ascent (descent)

basically go opposite the direction of the derivative.

Consider the objective function:  $J(x) = x^2 - 6x + 5$ 

gradient = fprime(old\_x)
move = gradient \* step
current\_x = old\_x - move





### good step size







### too big step size







### too small step size







### **Example: Linear Regression**

$$\hat(y)=f_ heta(x)= heta^T x$$

### **Cost Function:**

$$J( heta) = rac{1}{2} \sum_{i=1}^m (f_ heta(x^{(i)} - y^{(i)})^2$$



### **Gradient Descent**

$$heta:= heta-\eta
abla_ heta J( heta)= heta-\eta\sum_{i=1}^m
abla J_i( heta)$$

### where $\eta$ is the learning rate.

### ENTIRE DATASET NEEDED

for i in range(n epochs): params\_grad = evaluate\_gradient(loss\_function, data, params) params = params - learning rate \* params grad`



### $(\theta)$

000

### Linear Regression: Gradient Descent

$$heta_j:= heta_j+lpha\sum_{i=1}^m(y^{(i)}-f_ heta(x^{(i)}))x_j^{(i)}$$





### Stochastic Gradient Descent

$$heta:= heta-lpha
abla_ heta J_i( heta)$$

### ONE POINT AT A TIME

```
for i in range(nb_epochs):
  np.random.shuffle(data)
  for example in data:
    params_grad = evaluate_gradient(loss_function, example, params)
    params = params - learning_rate * params_grad
```

Mini-Batch: do some at a time



LR, SGD: 
$$heta_j := heta_j + lpha(y^{(i)} - f_ heta(x$$

- the risk surface changes at each gradient calculation
- thus things are noisy
- cumulated risk is smoother, can be used to compare to SGD
- epochs are now the number of times you revisit the full dataset
- shuffle in-between to provide even more stochasticity



 $(i))x_{j}^{(i)}$ 





# Is in-sample Approximating out-of-sample?





### Hoeffding's inequality

population fraction  $\mu$ , sample drawn with replacement, fraction  $\nu$ :

$$P(|
u-\mu|>\epsilon)\leq 2e^{-2\epsilon^2N}$$

For hypothesis h, identify 1 with  $h(x_i) \neq f(x_i)$  at sample  $x_i$ . Then  $\mu, \nu$  are population/sample error rates. Then,

$$P(|R_{in}(h)-R_{out}(h)|>\epsilon)\leq 2e^{-2}$$



 $2\epsilon^2 N$ 

- Hoeffding inequality holds ONCE we have picked a hypothesis h, as we need it to label the 1 and 0s.
- But over the training set we one by one pick all the models in the hypothesis space
- best fit g is among the h in  $\mathcal{H}$ , g must be  $h_1$  OR  $h_2$  OR....Say **effectively** M such choices:

$$P(|R_{in}(g)-R_{out}(g)|\geq\epsilon)<=\sum_{h_i\in\mathcal{H}}P(|R_{in}(h_i)-R_{out}(h_i))$$



### $||\geq\epsilon)<=2\,M\,e^{-2\epsilon^2N}$

### Hoeffding, repharased:

Now let  $\delta = 2 M e^{-2\epsilon^2 N}$ .

### Then, with probability $1 - \delta$ :

$$R_{out} <= R_{in} + \sqrt{rac{1}{2N}ln(rac{2M}{\delta})}$$

For finite effective hypothesis set size M,  $R_{out} \sim R_{in}$  as N larger..



### Training vs Test

- training error approximates out-of-sample error slowly
- is test set just another sample like the training sample?
- key observation: test set is looking at only one hypothesis because the fitting is already done on the training set. So M=1for this sample!

$$R_{out} <= R_{in} + \sqrt{rac{1}{2N_{test}}ln(rac{2}{\delta})}$$



# Is this still a test set? Trouble:

- no discussion on the error bars on our error estimates
- "visually fitting" a value of  $d \implies$  contaminated test set.
- The moment we use it in the learning process, it is not a test set.



### Training vs Test

- the test set does not have an optimistic bias like the training set(thats why the larger effective M factor)
- once you start fitting for things like d on the test set, you cant call it a test set any more since we lose tight guarantee.
- test set has a cost of less data in the training set and must thus fit a less complex model.

