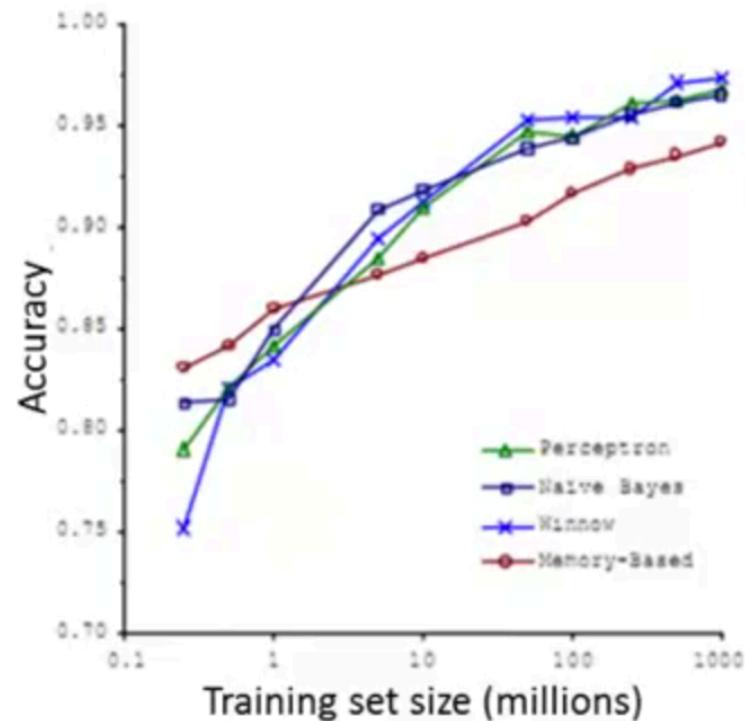


Machine learning and data

Classify between confusable words.
E.g., {to, two, too}, {then, than}.

For breakfast I ate _____ eggs.



“It’s not who has the best algorithm that wins.
It’s who has the most data.”

Linear regression with gradient descent

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

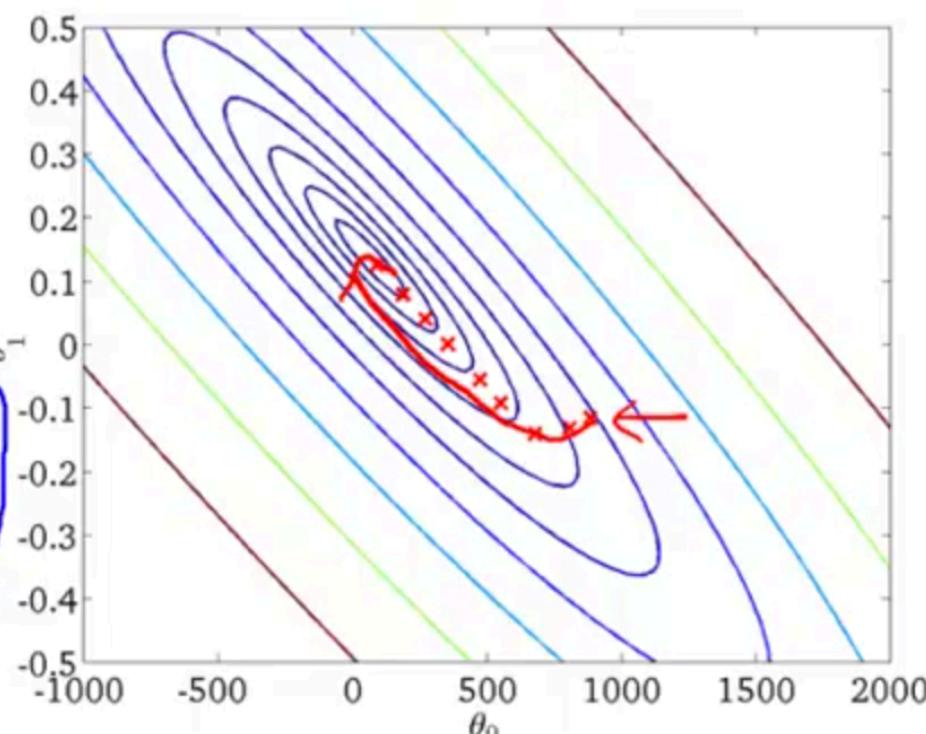
$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right]$$

(for every $j = 0, \dots, n$)

}

$$M = 300,000,000$$

Batch gradient descent



Batch gradient descent

$$\rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Repeat {

$$\rightarrow \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m \underbrace{(h_\theta(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_j} J_{train}(\theta)} x_j^{(i)}$$

(for every $j = 0, \dots, n$)

}

Stochastic gradient descent

$$\rightarrow \underbrace{cost(\theta, (x^{(i)}, y^{(i)}))}_{\uparrow} = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J_{train}(\theta) = \frac{1}{m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.

2. Repeat {

for $i=1, \dots, m$ {

$$\theta_j := \theta_j - \alpha \underbrace{(h_\theta(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_j} cost(\theta, (x^{(i)}, y^{(i)}))} \cdot x_j^{(i)}$$

} (for $j=0, \dots, n$)

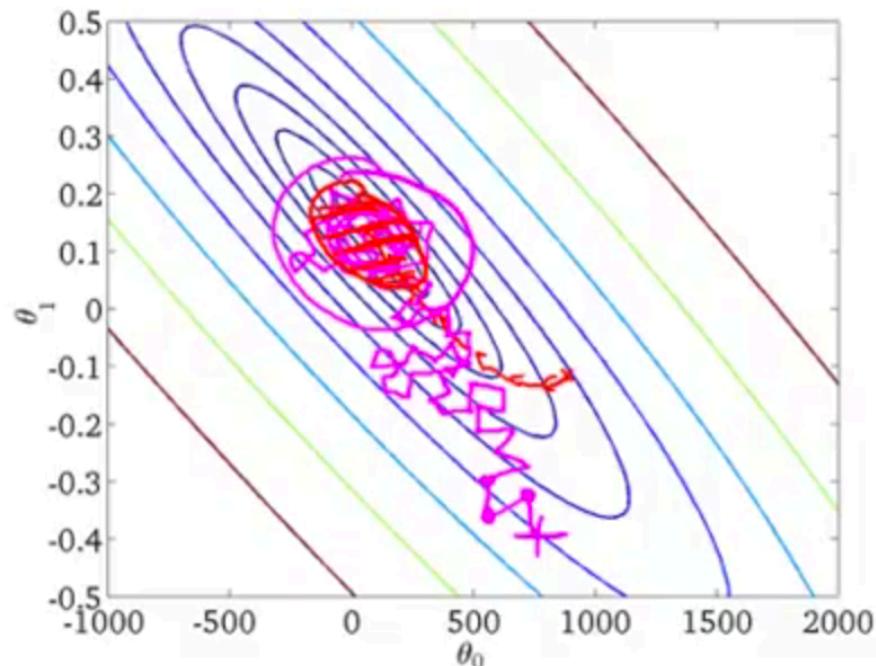
$$\rightarrow \frac{\partial}{\partial \theta_j} cost(\theta, (x^{(i)}, y^{(i)}))$$

Stochastic gradient descent

→ 1. Randomly shuffle (reorder) training examples

→ 2. Repeat {
 for $i := 1, \dots, m$ {
 $\theta_j := \theta_j - \alpha(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$
 (for every $j = 0, \dots, n$)
 }
}

$$\rightarrow m = 300,000,000$$



Mini-batch gradient descent

Say $b = 10$, $m = 1000$.

→ b examples

→ 1 example

Vectorization

Repeat {

→ for $i = 1, 11, 21, 31, \dots, 991$ {

→ $\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$

(for every $j = 0, \dots, n$)

}

}

$$m = \underline{300,000,000}$$



$$\frac{b=10}{\uparrow}$$

Checking for convergence

→ Batch gradient descent:

→ Plot $J_{train}(\theta)$ as a function of the number of iterations of gradient descent.

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

M = 300, 500, 1000

→ Stochastic gradient descent:

$$cost(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

During learning, compute $\frac{cost(\theta, (x^{(i)}, y^{(i)}))}{\uparrow}$ before updating θ using $(x^{(i)}, y^{(i)})$. \uparrow

$\Rightarrow (x^{(i)}, y^{(i)}) , (x^{(i+1)}, y^{(i+1)}) , \dots$

→ Every 1000 iterations (say), plot $cost(\theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1000 examples processed by algorithm.

Stochastic gradient descent

$$cost(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2}(h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.

2. Repeat {

```
for i := 1, ..., m      {  
    theta_j := theta_j - alpha(h_theta(x^{(i)}) - y^{(i)})x_j^{(i)}  
    (for j = 0, ..., n)  
}
```

```
}
```

Learning rate α is typically held constant. Can slowly decrease α over time if we want θ to converge. (E.g. $\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}$)

