

Week 9

Density Estimation

Problem Motivation

Anomaly detection example

- Fraud detection:
 - $x^{(i)}$ = features of user i 's activities
 - Model $p(x)$ from data.
 - Identify unusual users by checking which have $\underline{p(x) < \epsilon}$
- Manufacturing
- Monitoring computers in a data center.
 - $x^{(i)}$ = features of machine i
 - x_1 = memory use, x_2 = number of disk accesses/sec,
 x_3 = CPU load, x_4 = CPU load/network traffic.
 - ...

x_1
 x_2
 x_3
 x_4
 $p(x)$

Your anomaly detection system flags x as anomalous whenever $p(x) \leq \epsilon$. Suppose your system is flagging too many things as anomalous that are not actually so (similar to supervised learning, these mistakes are called false positives). What should you do?

- Try increasing ϵ .
- Try decreasing ϵ .

Correct

Gaussian Distribution

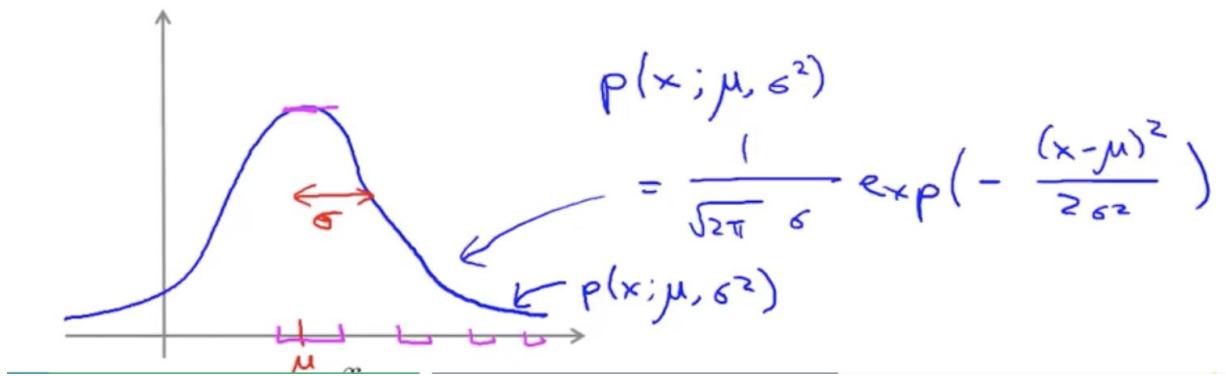
Gaussian (Normal) distribution

Say $x \in \mathbb{R}$. If x is a distributed Gaussian with mean μ , variance σ^2 .

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

"distributed as"

σ standard deviation



$$\underline{\mu} = \frac{1}{m} \sum_{i=1}^m \underline{x}^{(i)}$$

$$\underline{\sigma^2} = \frac{1}{m} \sum_{i=1}^m (\underline{x}^{(i)} - \underline{\mu})^2$$

Algorithm

Anomaly detection algorithm

- 1. Choose features x_i that you think might be indicative of anomalous examples. $\{x^{(1)}, \dots, x^{(m)}\}$
- 2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$
 - $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$ $p(x_j; \mu_j, \sigma_j^2)$
 - $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$
- 3. Given new example x , compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if $p(x) < \varepsilon$

Building an Anomaly Detection System

Developing and Evaluating an Anomaly Detection System

Aircraft engines motivating example

- > 10000 good (normal) engines
- > 20 flawed engines (anomalous) 2-50 y=1
- > Training set: 6000 good engines ($y=0$) $p(x) = p(x_1; \mu_1, \sigma^2_1) \dots p(x_n; \mu_n, \sigma^2_n)$
 CV: 2000 good engines ($y=0$), 10 anomalous ($y=1$)
 Test: 2000 good engines ($y=0$), 10 anomalous ($y=1$)

Algorithm evaluation

- > Fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(m)}\}$
- > On a cross validation/test example x , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (normal)} \end{cases}$$

Possible evaluation metrics:

- True positive, false positive, false negative, true negative
- Precision/Recall
- F_1 -score

Can also use cross validation set to choose parameter ε

$y=0$ have high accuracy, thus put the equal sign in the normal line

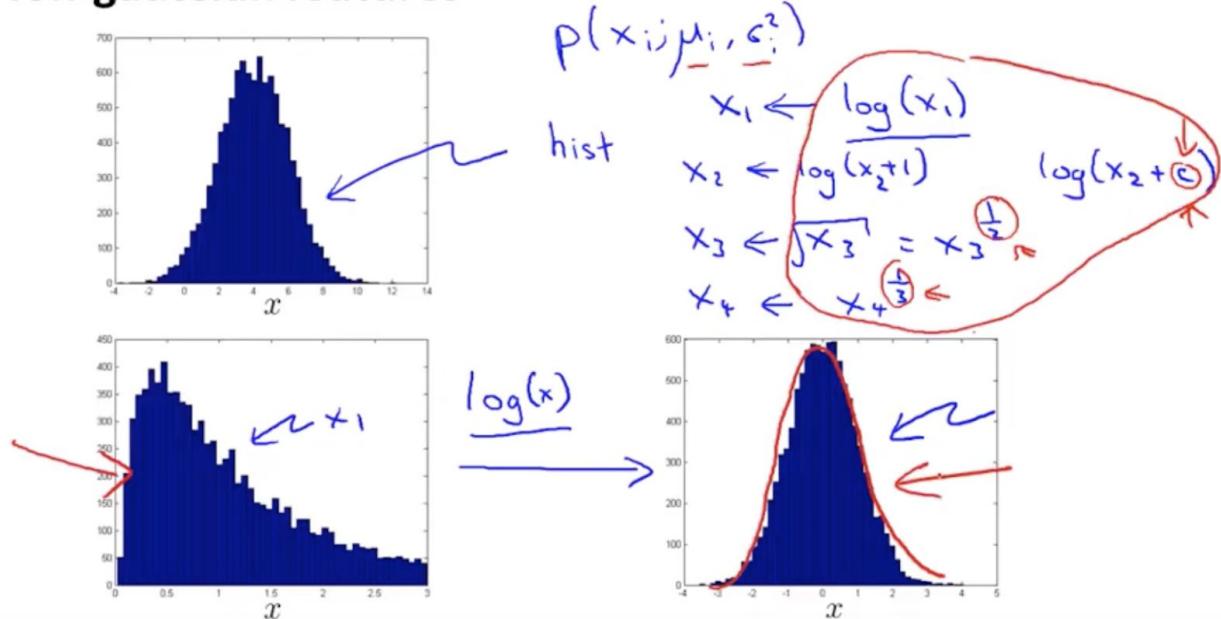
Anomaly Detection vs. Supervised Learning

Anomaly detection	vs.	Supervised learning
<p>Very small number of positive examples ($y = 1$). (0-20 is common).</p> <p>Large number of negative ($y = 0$) examples. $p(x)$</p> <p>Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we’ve seen so far.</p>		<p>Large number of positive and negative examples.</p> <p>Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.</p>

Anomaly detection	vs.	Supervised learning
<ul style="list-style-type: none"> Fraud detection $y=1$ Manufacturing (e.g. aircraft engines) Monitoring machines in a data center 		<ul style="list-style-type: none"> Email spam classification Weather prediction (sunny/rainy/etc). Cancer classification

Choosing What Features to Use

Non-gaussian features

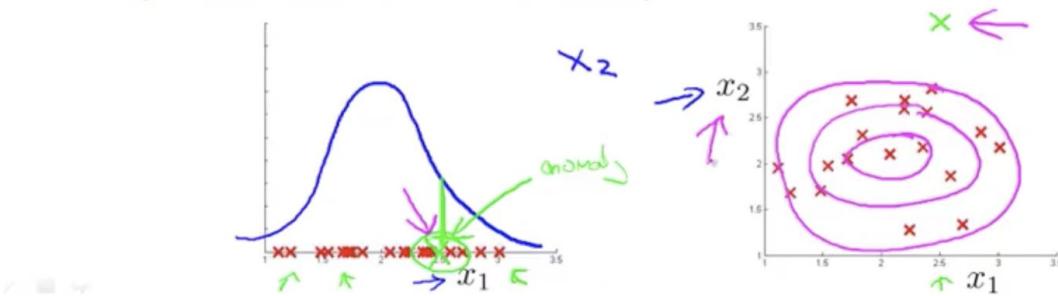


Error analysis for anomaly detection

- Want $p(x)$ large for normal examples x .
- $p(x)$ small for anomalous examples x .

Most common problem:

- $p(x)$ is comparable (say, both large) for normal and anomalous examples



Multivariate Gaussian Distribution

Multivariate Gaussian Distribution

Multivariate Gaussian (Normal) distribution

$x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2), \dots$, etc. separately.

Model $p(x)$ all in one go.

Parameters: $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

$$p(x; \mu, \Sigma) =$$

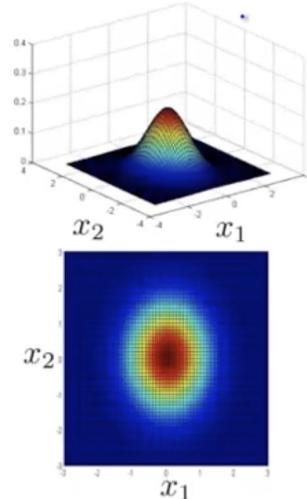
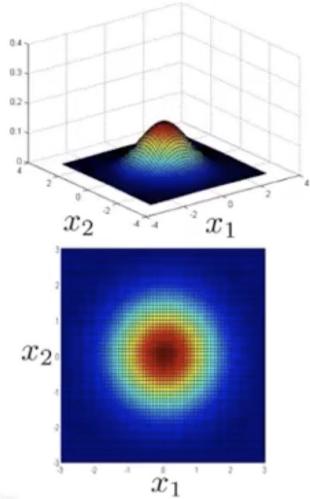
$$\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu)\right)$$

$| \Sigma | = \text{determinant of } \Sigma \quad | \det(\text{Sigma}) |$

Multivariate Gaussian (Normal) examples

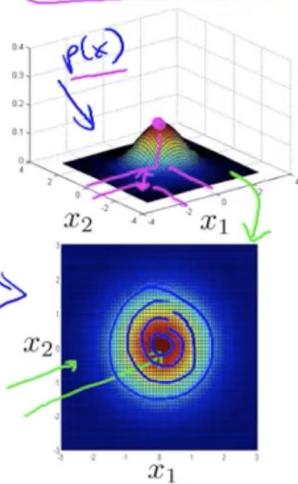
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$

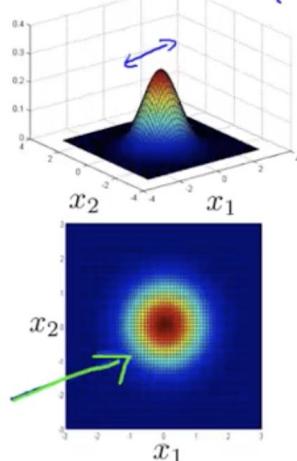


Multivariate Gaussian (Normal) examples

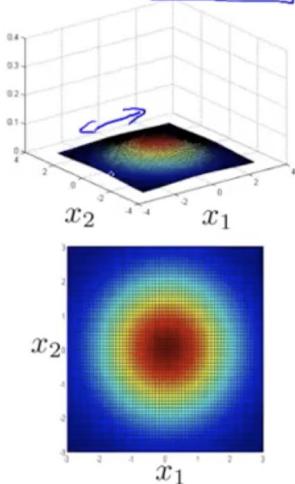
→ $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

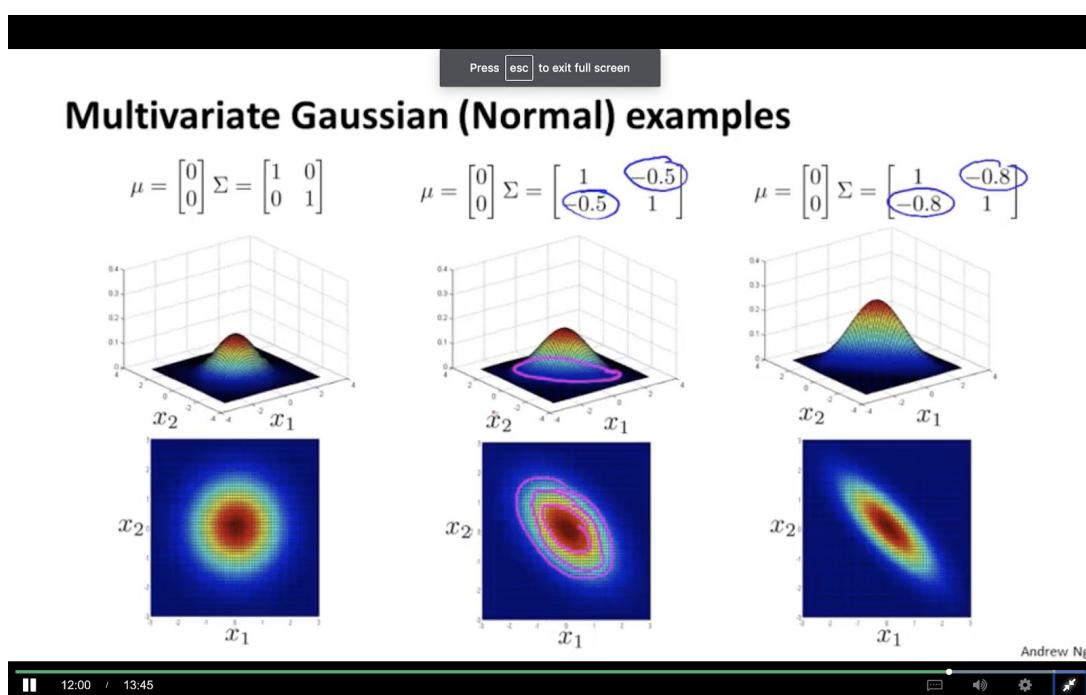


$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$





Anomaly Detection using the Multivariate Gaussian Distribution

Anomaly detection with the multivariate Gaussian

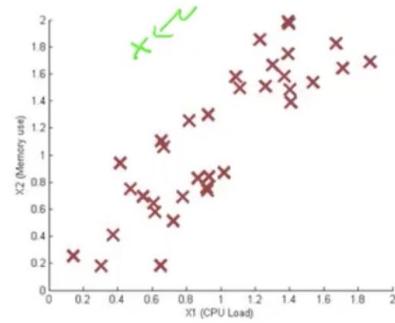
1. Fit model $p(x)$ by setting

$$\begin{cases} \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \\ \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T \end{cases}$$

2. Given a new example x , compute

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

Flag an anomaly if $\underline{p(x) < \varepsilon}$



Andrew Ng



$$\Sigma = \begin{bmatrix} \sigma_1^2 & \dots & \dots \\ \vdots & \ddots & \vdots \\ \dots & \dots & \sigma_n^2 \end{bmatrix}$$

⇒ Original model	vs.	⇒ Multivariate Gaussian
$p(x_1; \mu_1, \sigma_1^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$		$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} \Sigma ^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$
Manually create features to capture anomalies where x_1, x_2 take unusual combinations of values. → $X_3 = \frac{x_1}{x_2} = \frac{\text{CPU load}}{\text{memory}}$		→ Automatically captures correlations between features $\Sigma \in \mathbb{R}^{n \times n}$ Σ^{-1}
→ Computationally cheaper (alternatively, scales better to large n) $n=10,000, m=100,000$		Computationally more expensive
OK even if m (training set size) is small		Must have $m > n$, or else Σ is non-invertible.

Andrew Ng



The original model $p(x_1; \mu_1, \sigma_1^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$ corresponds to a multivariate Gaussian where the contours of $p(x; \mu, \Sigma)$ are axis-aligned.

Quiz

1. If you do not have any labeled data (or if all your data has label $y = 0$), then it is still possible to learn $p(x)$, but it may be harder to evaluate the system or choose a good value of ε
2. Anomaly detection only models the negative examples, whereas an SVM learns to discriminate between positive and negative examples, so the SVM will perform better when you have many positive and negative examples.

Predicting Movie Ratings

Problem Formulation

Recommender System

- n_u = no. users
- n_m = no. movies
- $r(i, j) = 1$ if user j has
rated movie i
- $y^{(i,j)}$ = rating given by
user j to movie i
(defined only if
 $r(i, j) = 1$)

Content Based Recommendations

Content-based recommender systems

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$n_u = 4, n_m = 5$	$x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$
Love at last 1	5	5	0	0	x_1 (romance) $\rightarrow 0.9$	x_2 (action) $\rightarrow 0$
Romance forever 2	5	?	?	0	$\rightarrow 1.0$	$\rightarrow 0.01$
Cute puppies of love 3	?	4	0	?	$\rightarrow 0.99$	$\rightarrow 0$
Nonstop car chases 4	0	0	5	4	$\rightarrow 0.1$	$\rightarrow 1.0$
Swords vs. karate 5	0	0	5	?	$\rightarrow 0$	$\rightarrow 0.9$

$n=2$

For each user j , learn a parameter $\theta^{(j)} \in \mathbb{R}^3$. Predict user j as rating movie i with $(\theta^{(j)})^T x^{(i)}$ stars.

$$\theta^{(j)} \in \mathbb{R}^{n+1}$$

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \Leftrightarrow \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad (\theta^{(1)})^T x^{(3)} = 5 \times 0.99 = 4.95$$

Andrew Ng

Problem formulation

- $r(i, j) = 1$ if user j has rated movie i (0 otherwise)
- $y^{(i,j)}$ = rating by user j on movie i (if defined)
- $\theta^{(j)}$ = parameter vector for user j
- $x^{(i)}$ = feature vector for movie i
- For user j , movie i , predicted rating: $\underline{(\theta^{(j)})^T(x^{(i)})}$ $\theta^{(j)} \in \mathbb{R}^{n+1}$
- $m^{(j)}$ = no. of movies rated by user j

To learn $\underline{\theta^{(j)}}$:

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: r(i,j)=1} ((\theta^{(j)})^T(x^{(i)}) - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (\theta^{(j)})_k^2$$

Andrew Ng

Optimization algorithm:

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

↙

$J(\theta^{(1)}, \dots, \theta^{(n_u)})$

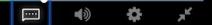
Gradient descent update:

$$\begin{aligned} \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k = 0) \\ \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0) \end{aligned}$$

$\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(1)}, \dots, \theta^{(n_u)})$

Andrew Ng

12:38 / 14:31



Collaborative Filtering

Collaborative Filtering

Recommender System

Collaborative filtering

Given $x^{(1)}, \dots, x^{(n_m)}$ (and movie ratings),
can estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$,
can estimate $x^{(1)}, \dots, x^{(n_m)}$

Collaborative Filtering Algorithm

Minimizing $x^{(1)}, \dots, x^{(n_m)}$ and $\theta^{(1)}, \dots, \theta^{(n_u)}$ simultaneously:

$$J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$$\min_{\substack{x^{(1)}, \dots, x^{(n_m)} \\ \theta^{(1)}, \dots, \theta^{(n_u)}}} J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$$

Andrew Ng

Collaborative filtering algorithm

- 1. Initialize $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ to small random values.
- 2. Minimize $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ using gradient descent (or an advanced optimization algorithm). E.g. for every $j = 1, \dots, n_u, i = 1, \dots, n_m$:

$$x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right)$$
- 3. For a user with parameters θ and a movie with (learned) features x , predict a star rating of $\theta^T x$.

In the algorithm we described, we initialized $x^{(1)}, \dots, x^{(n_m)}$ and $\theta^{(1)}, \dots, \theta^{(n_u)}$ to small random values. Why is this?

- This step is optional. Initializing to all 0's would work just as well.
- Random initialization is always necessary when using gradient descent on any problem.
- This ensures that $x^{(i)} \neq \theta^{(j)}$ for any i, j .
- This serves as symmetry breaking (similar to the random initialization of a neural network's parameters) and ensures the algorithm learns features $x^{(1)}, \dots, x^{(n_m)}$ that are different from each other.

Correct

Low Rank Matrix Factorization

Vectorization: Low Rank Matrix Factorization

Collaborative filtering

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

Predicted ratings: $(\Theta^{(i)})^T(x^{(j)})$

$$\times \Theta^T \quad (i,j) \uparrow$$

$$\left[\begin{array}{cccc} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \dots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \dots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \dots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{array} \right]$$

$\rightarrow X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(n_m)})^T \end{bmatrix} \quad \Theta = \begin{bmatrix} -(\Theta^{(1)})^T \\ -(\Theta^{(2)})^T \\ \vdots \\ -(\Theta^{(n_u)})^T \end{bmatrix}$

Low rank matrix factorization

Finding related movies

For each product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$.

$\rightarrow x_1 = \text{romance}, x_2 = \text{action}, x_3 = \text{comedy}, x_4 = \dots$

How to find $\underset{\uparrow}{\text{movies } j}$ related to $\underset{\uparrow}{\text{movie } i}$?

small $\|x^{(i)} - x^{(j)}\| \rightarrow$ movie j and i are "similar"

5 most similar movies to movie i :

Find the 5 movies j with the smallest $\|x^{(i)} - x^{(j)}\|$.

Implementational Detail: Mean Normalization

Finding related movies

For each product i , we learn a feature vector $\underline{x^{(i)}} \in \mathbb{R}^n$.

$\rightarrow x_1 = \text{romance}$, $x_2 = \text{action}$, $x_3 = \text{comedy}$, $x_4 = \dots$

How to find $\underline{\text{movies } j}$ related to $\underline{\text{movie } i}$?

small $\|x^{(i)} - x^{(j)}\|$ \rightarrow movie j and i are "similar"

5 most similar movies to movie i :

Find the 5 movies j with the smallest $\|x^{(i)} - x^{(j)}\|$.

Quiz

Suppose you run a bookstore, and have ratings (1 to 5 stars) of books. Your collaborative filtering algorithm has learned a parameter vector $\theta^{(j)}$ for user j , and a feature vector $x^{(i)}$ for each book. You would like to compute the "training error", meaning the average squared error of your system's predictions on all the ratings that you have gotten from your users. Which of these are correct ways of doing so (check all that apply)?

For this problem, let m be the total number of ratings you have gotten from your users. (Another way of saying this is

that $m = \sum_{i=1}^{n_m} \sum_{j=1}^{n_u} r(i, j)$). [Hint: Two of the four options below are correct.]

- $\frac{1}{m} \sum_{(i,j):r(i,j)=1} (\sum_{k=1}^n (\theta^{(j)})_k x_k^{(i)} - y^{(i,j)})^2$
- $\frac{1}{m} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} (\sum_{k=1}^n (\theta^{(j)})_k x_k^{(i)} - y^{(i,j)})^2$

Suppose you have two matrices A and B , where A is 5×3 and B is 3×5 . Their product is $C = AB$, a 5×5 matrix. Furthermore, you have a 5×5 matrix R where every entry is 0 or 1. You want to find the sum of all elements $C(i, j)$ for which the corresponding $R(i, j)$ is 1, and ignore all elements $C(i, j)$ where $R(i, j) = 0$. One way to do so is the following code:

```
C = A * B;
total = 0;
for i = 1:5
    for j = 1:5
        if (R(i,j) == 1)
            total = total + C(i,j);
        end
    end
end
```

Which of the following pieces of Octave code will also correctly compute this total? Check all that apply. Assume all options are in code.

- A. `total = sum(sum((A * B) .* R))`
- B. `C = A * B; total = sum(sum(C(R == 1)));`



Week 10

Gradient Descent with Large Datasets

Learning With Large Datasets

Stochastic Gradient Descent

Linear regression with gradient descent

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

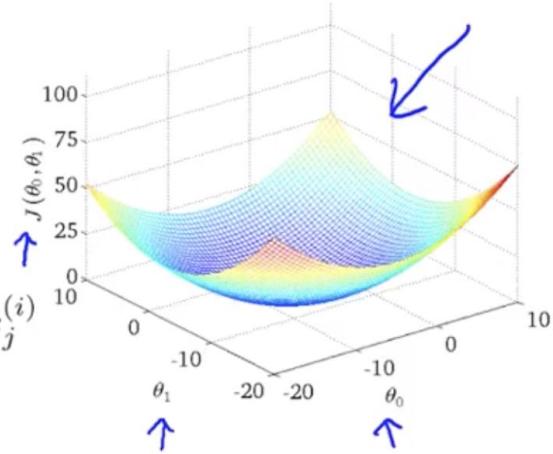
$$\Rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

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

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

}



Batch gradient descent

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

Repeat {

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

$\frac{\partial}{\partial \theta_j} J_{train}(\theta)$

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

}

Stochastic gradient descent

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

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

1. Randomly shuffle dataset. ←

2. Repeat {

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

$$\theta_j := \theta_j - \alpha \frac{(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}}{(for j=0, \dots, n)}$$

}

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

$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots$

Mini-Batch Gradient Descent

Mini-batch gradient descent

- Batch gradient descent: Use all m examples in each iteration
- Stochastic gradient descent: Use 1 example in each iteration

Mini-batch gradient descent: Use b examples in each iteration

$$\begin{aligned} b &= \text{mini-batch size.} & b &= 10. & 2 - 100 \\ \text{Get } b = 10 \text{ examples} & & (x^{(i)}, y^{(i)}), \dots, (x^{(i+9)}, y^{(i+9)}) & & \\ \rightarrow \theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) \cdot x_j^{(k)} & & & & \\ j := i + 10 & & & & \end{aligned}$$

Stochastic Gradient Descent Convergence

Checking for convergence

Batch gradient descent:

- Plot $J_{train}(\theta)$ as a function of the number of iterations of gradient descent.
- $$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

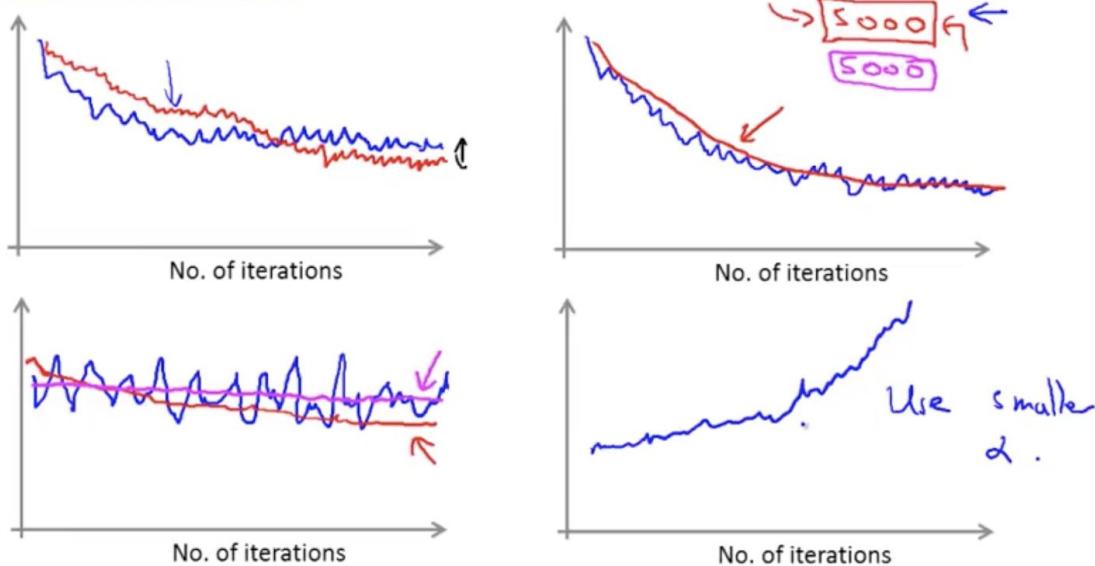
$M = 300, 500, 500$

Stochastic gradient descent:

- $$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$
- During learning, compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ using $(x^{(i)}, y^{(i)})$.
- Every 1000 iterations (say), plot $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1000 examples processed by algorithm.

Checking for convergence

Plot $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$, averaged over the last 1000 (say) examples



Ans

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

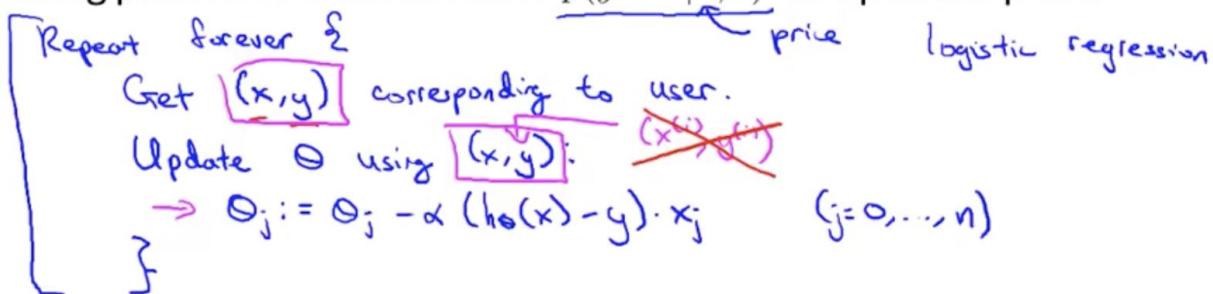
Advanced Topics

Online Learning

Online learning

Shipping service website where user comes, specifies origin and destination, you offer to ship their package for some asking price, and users sometimes choose to use your shipping service ($y = 1$), sometimes not ($y = 0$).

Features x capture properties of user, of origin/destination and asking price. We want to learn $p(y = 1|x; \theta)$ to optimize price.



Can adopt to changing user preferences.

Other online learning example:

Product search (learning to search)

User searches for "Android phone 1080p camera" \leftarrow

Have 100 phones in store. Will return 10 results.

$\rightarrow x = \text{features of phone, how many words in user query match name of phone, how many words in query match description of phone, etc.}$

$\rightarrow y = 1$ if user clicks on link. $y = 0$ otherwise.

\rightarrow Learn $p(y = 1|x; \theta)$. \leftarrow predicted CTR

\rightarrow Use to show user the 10 phones they're most likely to click on.

$(x, y) \leftarrow$
↑ ↑

CTR: click through rate

Other examples: Choosing special offers to show user; customized selection of news articles; product recommendation; ...

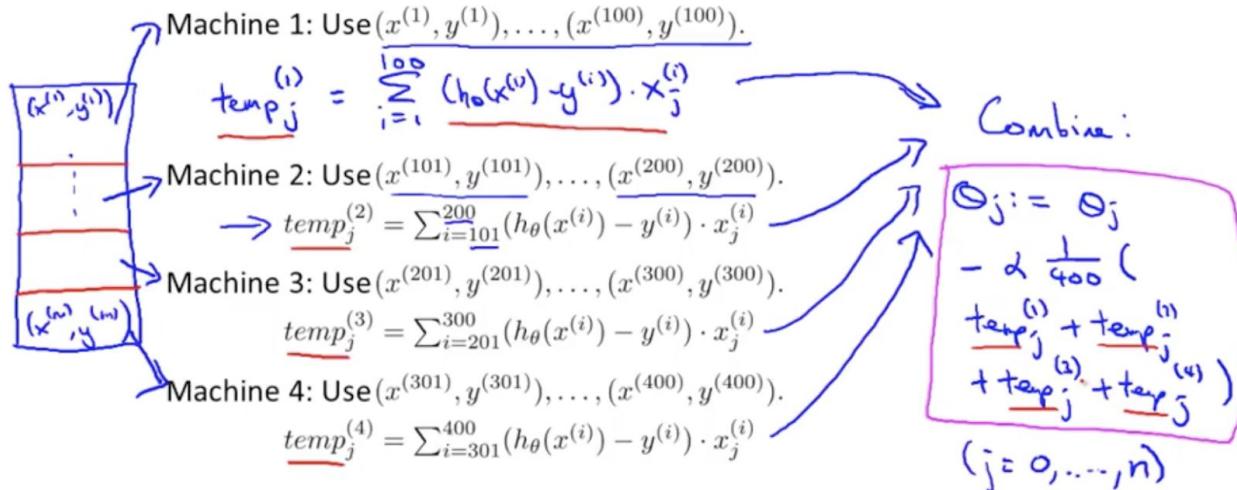
Map Reduce and Data Parallelism

Map-reduceBatch gradient descent:

$$m = 400$$

$$m = 400,000,000$$

$$\theta_j := \theta_j - \alpha \frac{1}{400} \sum_{i=1}^{400} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

**Map-reduce and summation over the training set**

Many learning algorithms can be expressed as computing sums of functions over the training set.

E.g. for advanced optimization, with logistic regression, need:

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

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

\nwarrow

$\text{temp}_j^{(1)} \quad \text{temp}_j^{(1)} \leftarrow$

Suppose you apply the map-reduce method to train a neural network on ten machines. In each iteration, what will each of the machines do?

- Compute either forward propagation or back propagation on 1/5 of the data.
- Compute forward propagation and back propagation on 1/10 of the data to compute the derivative with respect to that 1/10 of the data.

Correct

Quiz

4. Assuming that you have a very large training set, which of the following algorithms do you think can be parallelized using map-reduce and splitting the training set across different machines? Check all that apply.

- CORRECT A neural network trained using batch gradient descent.
- WRONG Logistic regression trained using stochastic gradient descent.
- CORRECT Linear regression trained using batch gradient descent.
- WRONG An online learning setting, where you repeatedly get a single example (x,y) , and want to learn from that single example before moving on.
- WRONG Computing the average of all the features in your training set

5. Which of the following statements about map-reduce are true? Check all that apply.

- WRONG Running map-reduce over N computers requires that we split the training set into N^2 pieces.
- **right** WRONG In order to parallelize a learning algorithm using map-reduce, the first step is to figure out how to express the main work done by the algorithm as computing sums of functions of training examples.
- CORRECT When using map-reduce with gradient descent, we usually use a single machine that accumulates the gradients from each of the map-reduce machines, in order to compute the parameter update for that iteration.
- CORRECT If you are have just 1 computer, but your computer has multiple CPUs or multiple cores, then map-reduce might be a viable way to parallelize your learning algorithm.



Week 11

Application Example: Photo OCR

Problem Description and Pipeline

Photo OCR: Photo Optical Character Recognition

Pipeline: text detection, character segmentation, character classification

Sliding Windows

Getting Lots of Data and Artificial Data

Discussion on getting more data

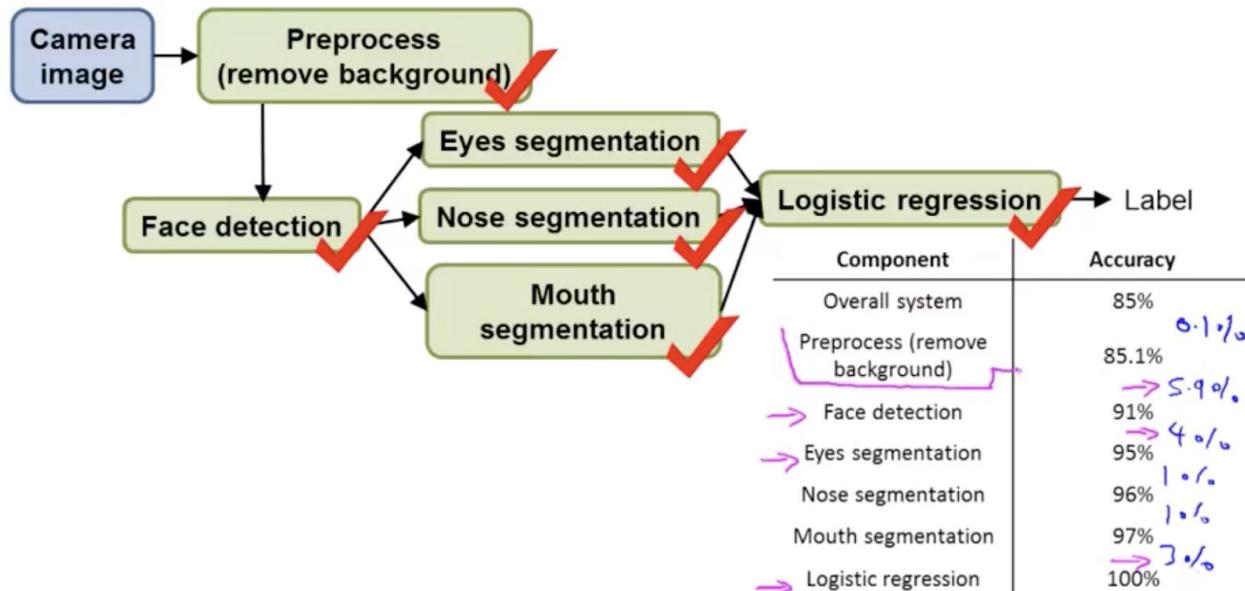
1. Make sure you have a low bias classifier before expending the effort. (Plot learning curves). E.g. keep increasing the number of features/number of hidden units in neural network until you have a low bias classifier.
2. “How much work would it be to get 10x as much data as we currently have?”
 - Artificial data synthesis
 - Collect/label it yourself
 - “Crowd source” (E.g. Amazon Mechanical Turk)

Ceiling Analysis: What Part of the Pipeline to Work on Next

Component	Accuracy
Overall system	72%
→ Text detection	89% ↕ 17%
Character segmentation	90% ↕ 1%
Character recognition	100% ↕ 10%

Text detection对整个系统有17%的提升, character segmentation(1%), character recognition(10%)

Another ceiling analysis example



Summary

Summary: Main topics

› Supervised Learning

- Linear regression, logistic regression, neural networks, SVMs

$$(x^{(i)}, y^{(i)})$$

Unsupervised Learning

- K-means, PCA, Anomaly detection

$$x^{(i)}$$

Special applications/special topics

- Recommender systems, large scale machine learning.

› Advice on building a machine learning system

- Bias/variance, regularization; deciding what to work on next: evaluation of learning algorithms, learning curves, error analysis, ceiling analysis.

