

Week 6

Advice for Applying ML

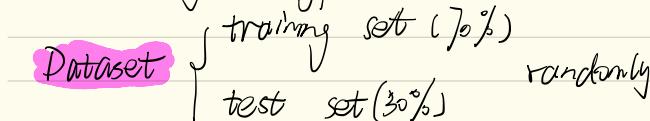
1. Evaluating a Learning Algorithm

1. Deciding what to try next

- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features (x_1^2, x_2^2, x_1x_2 , etc.)
- Try decreasing λ
- Try increasing λ

$$x_1, x_2, x_3, \dots, x_{100}$$

2. Evaluating a Hypothesis



Steps.

<> Learn θ and minimize $J(\theta)$ \Leftarrow training emps.

<> Compute the test set error $J_{\text{test}}(\theta)$.

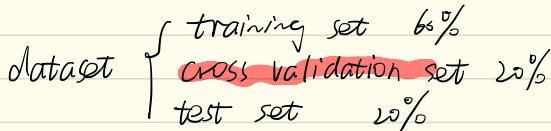
$J_{\text{test}}(\theta)$

<> Linear: $J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x^{(i)}_{\text{test}}) - y^{(i)}_{\text{test}})^2$

<> Misclassification: (0, 1) 类型.

$\text{err}(h_{\theta}(x), y) = \begin{cases} 1, & h_{\theta}(x) > 0.5 \text{ and } y = 0, \text{ or } h_{\theta}(x) < 0.5 \text{ and } y = 1 \\ 0, & \text{otherwise} \end{cases}$

3. Model Selection and Train / Validation / Test sets.



1. Optimize the parameters in Θ using the training set for each polynomial degree.

2. Find the polynomial degree d with the least error using the cross validation set.

3. Estimate the generalization error using the test set with $J_{\text{test}}(\Theta^{(d)})$, ($d = \theta$ from polynomial with lower error);

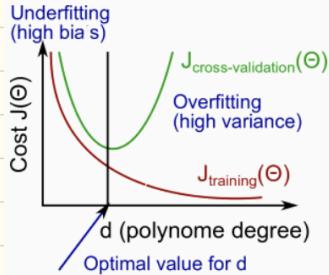
Week 6

Advice for Applying ML

2. Bias vs. Variance

1. Diagnosing Bias vs. Variance

→ 模型不夠的起由 bias 及 variance 何者？



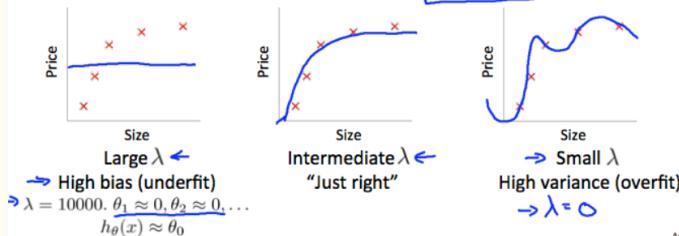
↓ High bias (underfitting)
↓ High variance (overfitting)

2. Regularization and Bias / Variance

Linear regression with regularization

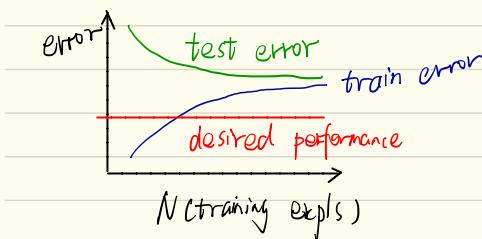
$$\text{Model: } h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \leftarrow$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2 \quad \leftarrow$$



3. Learning curves

① High bias



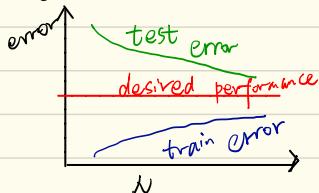
training set size ↓ : $J_{train}(\theta)$ ↓
 $J_{cv}(\theta)$ ↗

training set size ↑ : $J_{train}(\theta)$ ↑
 $J_{cv}(\theta)$ ↑

Week 6

Advice for Applying ML

② High Variance



training set size N : $J_{\text{train}}(\theta)$ 小 $J_{\text{CV}}(\theta)$ 大
training set size T : $J_{\text{train}}(\theta)$ 增大 $J_{\text{CV}}(\theta)$ 减小
 $J_{\text{train}}(\theta) < J_{\text{CV}}(\theta)$
且相差很大

④ Deciding What to Do Next Revised

- **Getting more training examples:** Fixes high variance
- **Trying smaller sets of features:** Fixes high variance
- **Adding features:** Fixes high bias
- **Adding polynomial features:** Fixes high bias
- **Decreasing λ :** Fixes high bias
- **Increasing λ :** Fixes high variance.

⑤ neutral network.

左 θ : underfitting high bias

右 θ : overfitting high variance

TF网 θ 层 θ 单元

Model Complexity Effects.

- Lower-order polynomials (low model complexity) have high bias and low variance. In this case, the model fits poorly consistently.
- Higher-order polynomials (high model complexity) fit the training data extremely well and the test data extremely poorly. These have low bias on the training data, but very high variance.
- In reality, we would want to choose a model somewhere in between, that can generalize well but also fits the data reasonably well.

三. Building a Spam Classifier

Week 6

Advice for Applying ML

1. Prioritizing What to Work On

16万 / 10000 ~ 50000 篇 出现 (email) word

→ X feature

| A lot of **data**

| Sophisticated **features**

| **Algorithms**

2. Error Analysis

首先找一个极简模型，初步实现。

以决定 where to cost time next

Recommended approach to solving ML problem

stemming 特殊行进法

- Start with a simple algorithm, implement it quickly, and test it early on your cross validation data.
- Plot learning curves to decide if more data, more features, etc. are likely to help.
- Manually examine the errors on examples in the cross validation set and try to spot a trend where most of the errors were made.

3. Handling Skewed Data

1. Error Metrics for Skewed Classes

		Actual class	
		1	0
Predicted class	1	True Positive False Positive	False Negative
	0	True Negative	False Neg

$$\text{precision} = \frac{\text{True Pos}}{\#\text{predicted pos}} = \frac{\text{True Pos}}{\text{True Pos} + \text{False Pos}}$$

$$\text{Recall} = \frac{\text{True Pos}}{\#\text{actual Pos}} = \frac{\text{True Pos}}{\text{True Pos} + \text{False Neg}}$$

2. Trading Off precision and recall

Week 6

Advice for Applying ML

logistic regression classifier

1) Predict $y=1$ if $h_\theta(x) \geq \text{threshold}$

2) Predict $y=0$ if $h_\theta(x) < \text{threshold}$

way to determine threshold.

$$F_1 \text{ Score} = 2 \frac{PR}{P+R} \quad (\text{maximize } F_1 \text{ Score})$$

i.e. P (Precision)

R (recall) from cross validation set

IV. Using Large Data Sets

It's not who has the best algorithm that wins,
but who has the most training sets

Week 7

Support Vector Machines

1. Large Margin Classification

1. Optimization Objective

Logistic regression:

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \underbrace{\left(-\log h_{\theta}(x^{(i)}) \right)}_{\text{cost}_1(\theta^T x^{(i)})} + (1 - y^{(i)}) \underbrace{\left(-\log(1 - h_{\theta}(x^{(i)})) \right)}_{\text{cost}_0(\theta^T x^{(i)})} \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Support vector machine

$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^n \theta_j^2$$

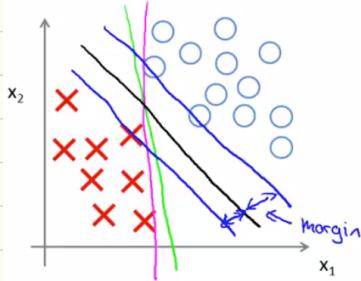
↓
 $C = \frac{1}{\lambda}$

SVM Hypothesis:

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

2. Large Margin Intuition

SVM Decision Boundary: Linearly separable case



Large margin classifier

2. Kernels

Compute new features depending on landmarks.

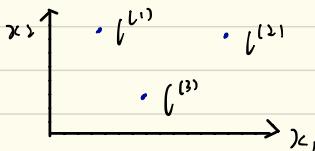
Week 7

Support Vector Machines

Kernels and Similarity

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{2\sigma^2}\right)$$

if $x \approx l^{(1)}$: $f_1 \approx 1$
if x far from $l^{(1)}$: $f_1 \approx 0$



SVM with Kernels

Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$,
choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$.

Given example x :

$$f_1 = \text{similarity}(x, l^{(1)})$$

$$f_2 = \text{similarity}(x, l^{(2)})$$

...

For training example $(x^{(i)}, y^{(i)})$:

$$f_i^{(1)} = \text{sim}(x^{(i)}, l^{(1)})$$

$$x^{(i)} : f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)})$$

$$\vdots \quad \quad \quad f_i^{(m)} = \text{sim}(x^{(i)}, l^{(m)}) = \exp(\dots).$$

$$f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)})$$

SVM Parameters

<1> $C (= \frac{1}{\lambda})$ \int Large C : Big variance & small bias
Small C : \propto λ

<2> σ^2 \int Large σ^2 : f_i vary more smoothly
 \Rightarrow Higher bias & lower variance
Small σ^2 : \propto λ

Week 7

Support Vector Machines

SVM in Practice

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters θ .

Need to specify:

→ Choice of parameter C.

Choice of kernel (similarity function):

E.g. No kernel ("linear kernel")

Predict "y = 1" if $\underline{\theta^T x} \geq 0$

$$\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0 \quad \rightarrow \underline{n \text{ large}}, \underline{m \text{ small}} \quad x \in \mathbb{R}^{n+1}$$

Gaussian kernel:

$$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right), \text{ where } l^{(i)} = x^{(i)}. \quad x \in \mathbb{R}^n, n \text{ small}$$

Need to choose $\underline{\sigma^2}$.

↑

and/or m large



Logistic regression vs. SVMs

n = number of features ($x \in \mathbb{R}^{n+1}$), m = number of training examples

If n is large (relative to m): (e.g. $n \geq m$, $n = 10,000$, $m = 10 \dots 1000$)

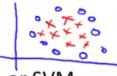
Use logistic regression, or SVM without a kernel ("linear kernel")

If n is small, m is intermediate: ($n = 1 \dots 1000$, $m = 10 \dots 10,000$) ←

→ Use SVM with Gaussian kernel

If n is small, m is large: ($n = 1 \dots 1000$, $m = 50,000 \dots$)

→ Create/add more features, then use SVM without a kernel



Neural network likely to work well for most of these settings, but may be slower to train.

Week 8

Unsupervised Learning

1. Clustering

unsupervised learning. input only x . no y .
cluster is an example of unsupervised ~.

1. k-means Algorithm

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

for $i = 1$ to m

$c^{(i)} :=$ index (from 1 to K) of cluster centroid
closest to $x^{(i)}$

for $k = 1$ to K

$\mu_k :=$ average (mean) of points assigned to cluster k

e.g. $k=2$ $x^{(1)}, x^{(3)}, x^{(5)}, x^{(6)} \Rightarrow c^{(1)}=2, c^{(3)}=2 \dots$

$$\mu_2 = \frac{1}{4} [x^{(1)}, x^{(3)}, x^{(5)}, x^{(6)}] \in \mathbb{R}^n$$

} cluster assignment

} move centroid

2. Optimization Objective

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2 \quad \leftarrow \min$$

Steps

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

for $i = 1$ to m

$c^{(i)} :=$ index (from 1 to K) of cluster centroid
closest to $x^{(i)}$

for $k = 1$ to K

$\mu_k :=$ average (mean) of points assigned to cluster k

3. Random Initialization

① Should have $K < m$

② Randomly pick K examples

③ Set μ_1, \dots, μ_K equal to these examples

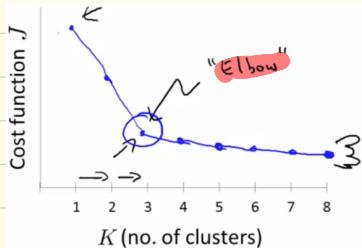
Week 8

Unsupervised Learning

4. Choosing the Number of Clusters

Choosing the value of K

Elbow Method



在 elbow 那里停止
就是最佳的。

试着把 T-shirt 从 5×1

变成 1×5

2. Motivation

1. Motivation I : Data Compression

Suppose we apply dimensionality reduction to a dataset of m examples $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$, where $x^{(i)} \in \mathbb{R}^n$. As a result of this, we will get out:

- A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(k)}\}$ of k examples where $k \leq n$.
- A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(k)}\}$ of k examples where $k > n$.
- A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$ of m examples where $z^{(i)} \in \mathbb{R}^k$ for some value of k and $k \leq n$.

Correct

- A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$ of m examples

2. Motivation II: Visualization

3. Principal Component Analysis

1. PCA Problem

Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error.

Week 8

Unsupervised Learning

PCA is Not Linear Regression

- ① L-r. x, y 方向 "投影".
 ② PCA. 垂直于 vector 的方向

2 - PCA Algorithm

PCA Algorithm Summary

After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T$$

$X = \begin{bmatrix} -x^{(1)^T} \\ \vdots \\ -x^{(m)^T} \end{bmatrix}$

$[U, S, V] = svd(\Sigma);$

$U_{reduce} = U(:, 1:k);$

$z = U_{reduce}' * x;$

$\uparrow \qquad \qquad \qquad x \in \mathbb{R}^n \quad \cancel{x \in \mathbb{R}}$

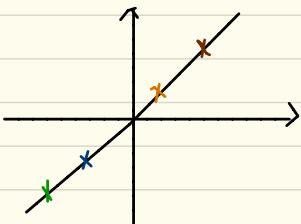
$$\Sigma = (1/m) * X' * X;$$

④ Applying PCA

1. Reconstruction from Compressed (representation)

$$\vec{z} = \underbrace{\sum_{\text{reduce}}}_{\text{compute reduction}} x$$

$$X_{\text{approx}} = U_{\text{reduce}} \cdot Z$$



2. Choosing the Number of Principal Components

$$\text{平均} \sum_{i=1}^m |x^{(i)} - x_{\text{approx}}^{(i)}|^2$$

k (dimensions)

$$\text{Total Variation} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

Week 8

Unsupervised Learning

Typically, choose k to be smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

$$\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2 \quad (0.05, 0.10 \text{ 部分}, \text{保留})$$

" 99% of variance is retained " (99% 保有)

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01 \quad (\text{For given } k)$$

`[U, S, V] = svd(Sigma)`

Pick smallest value of k for which

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

(99% of variance retained)

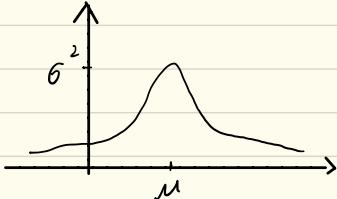
Week 9

Anomaly Detection

1. Density Estimation

1. Problem Motivation

2. Gaussian (Normal) Distribution (正态分布)



$$p(x) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Parameter Estimation

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

For given dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$

3. Algorithm

Anomaly detection algorithm

1. Choose features x_i that you think might be indicative of anomalous examples.
 $\{x^{(1)}, \dots, x^{(n)}\}$

2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\begin{aligned} \mu_j &= \frac{1}{m} \sum_{i=1}^m x_j^{(i)} & p(x_j; \mu_j, \sigma_j^2) &= \mu_1, \mu_2, \dots, \mu_n \\ \sigma_j^2 &= \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2 & \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} &= \frac{1}{m} \sum_{i=1}^m x^{(i)} \end{aligned}$$

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 $p(x) < \varepsilon$

2. Building an Anomaly Detection System.

training sets of ~~not~~ good (normal) examples

Week 9

Anomaly Detection

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
- Classification problem*

Can also use cross validation set to choose parameter ε

2. Anomaly Detection vs. Supervised Learning

Anomaly detection	vs.	Supervised learning
Very small number of positive examples ($y = 1$). (0-20 is common).		Large number of positive and negative examples. 
Large number of negative ($y = 0$) examples.  		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. 
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.		<i>Spam</i>  <ul style="list-style-type: none">• Email spam classification• Weather prediction (sunny/rainy/etc).• Cancer classification
Fraud detection 		
Manufacturing (e.g. aircraft engines)		
Monitoring machines in a data center		
	⋮	⋮

3. Choosing what futures to use

Week 9

Anomaly Detection

2. Predicting Movie Ratings

1. Content Based Recommendations

电影类型.

$$r(i, j) = 1, \text{ if user } j \text{ has rated movie } i \text{ (0 otherwise)}$$

$$y^{(i,j)} = \text{rating by user } j \text{ on movie } i \text{ (if defined)}$$

$\theta^{(j)}$ = parameter vector for user j .

$x^{(i)}$ = feature vector for movie i

For user j , movie i . predicted rating: $(\theta^{(j)})^T (x^{(i)})$

$m^{(j)}$ = no. of movies rated by user j .

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

Optimization objective:

To learn $\theta^{(j)}$ (parameter for user j):

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

To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}$:

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

Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k = 0)$$

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

3. Collaborative Filtering

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

Week 9

Anomaly Detection

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

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

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

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

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

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

Andrew Ng

协同过滤算法

<1> 简单地把 $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ 放在一起训练

<2> minimize $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ (梯度下降)

<3> 对 user, θ , movie, feature x .

使用 $\theta^T x$ 表示评分

IV. Low Rank Matrix Factorization

1. Vectorization

协同过滤：相似度矩阵法（基于用户，基于物品）

Predicted ratings:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \dots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \dots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \dots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{bmatrix}$$

2. Implementation Detail: Mean Normalization

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix} \rightarrow \mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

For user j , on movie i predict:

$$(\theta^{(i)})^T(x^{(i)}) + \mu_i$$

learn $\theta^{(i)}$

Week 10

Large Scale Machine Learning

1. Gradient Descent with Large Datasets

plot $J_{\text{train}}(\theta)$ $\Rightarrow J_{\text{val}}(\theta)$ \rightarrow 表明 m が大きいほど収束するまで時間がかかる

1. Stochastic Gradient Descent

m が大きい。

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

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

Steps:

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

2. 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

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

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

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

Week 10

Large Scale Machine Learning

3. Stochastic Gradient Descent Convergence

Check for convergence

Batch gradient descent:

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

$$\rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \quad M = 300, 000, 000$$

Stochastic gradient descent:

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

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

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

α is typically held constant. Can slowly decrease α over time if we want θ to converge.

(e.g. $\alpha = \frac{\text{const}}{\text{iterationNumber} + \text{const}}$)

2. Advanced Topics

1. Online Learning

Example :

Product Search / Recommendation. ;

Choosing special offers to user ;

Customized selection of news articles

2. Map-Reduce and Data Parallelism

分布式处理 · 分布式计算

Map-Reduce

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

Machine 1 : $1 \sim 100$

Machine 2 : $101 \sim 200$

:

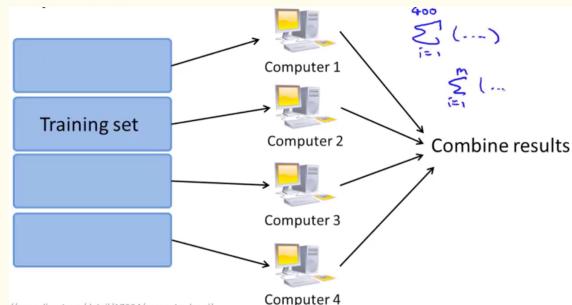
4 : $301 \sim 400$

$\left. \begin{array}{c} \\ \\ \\ \end{array} \right\}$ add together

$m = 400$

Week 10

Large Scale Machine Learning



https://www.csie.ntu.edu.tw/~cjlin/libSVM/3.20.3.4.LearningMachine.html#11

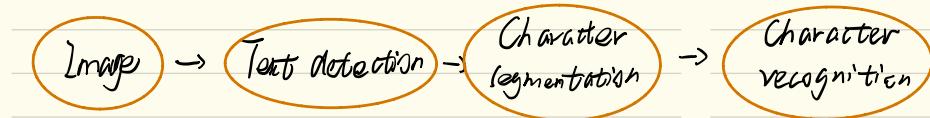
Idea: Many learning algorithm can be expressed as computing sum of functions over the training set.

Week 11

App Example: Photo OCR

- Photo OCR

1. Photo OCR Pipeline



2. Sliding Windows

1. Text detection



2. Character segmentation



3. Character classification



3. 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)

Summary

Main topics

Supervised Learning

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

- Linear regression, logistic regression, neural networks, SVMs

Unsupervised Learning

$$x^{(i)}$$

- K-means, PCA, Anomaly detection

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.