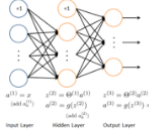


## Couse Summary

## Machine Learning Algorithm

ML Algorithm	Supervised Learning				Unsupervised Learning		
	Linear Regression	Logistic Regression	Neural Network	Support-Vector Machines (SVMs with Kernels)	K-Means	Principle Component Analysis (PCA)	Anomaly Detection
Motivation	Continuous output	Discrete output Classification problem (K-class)	Non-linear hypothesis w/o exhausting all non-linear terms	SVMs: Large margin classifier Kernel: non-linear decision boundary => linear	Most popular clustering algorithms	Dimensionality reduction for data compression/visualization	Very small # of positive examples
Input	$x_j^{(i)}$ - value j of i-th training example $X - m \times (n+1)$ matrix, row is each training example, and column is each feature (add bias feature 1 to first column)	$X - m \times (n+1)$ matrix (bias feature 1 added to first column)	$X - m \times (n+1)$ matrix, include bias unit Input layer (n units + 1 bias units), hidden layers (layer j has $S_j$ units)	$F - m \times (n+1)$ matrix (Kernels of X refer to landmarks L, initialization: L = X) (preprocessing X - $m \times n$ matrix -> F)	$X - m \times n$ matrix $M - K \times n$ matrix (random pick K training examples as <b>cluster centroids</b> )	$X - m \times n$ matrix	$X - m \times n$ matrix
Output	$y - m \times 1$ vector, row is each training example	$y - m \times K$ matrix, $y^{(i)} = [0 \ 0 \dots 1 \dots 0]$ , $y_k^{(i)} = 1$ ( <b>one-hot notation</b> ) (or $y - m \times 1$ vector, $y^{(i)} = k$ belongs to class k)	Assume a classification problem: $y - m \times K$ matrix, $y^{(i)} = [0 \ 0 \dots 1 \dots 0]$ output layer (K units)	$y - m \times K$ matrix, $y^{(i)} = [0 \ 0 \dots 1 \dots 0]$ , $y_k^{(i)} = 1$ (class k)	-- unsupervised -- Algo predicts the output as a $m \times 1$ vector, row is each training example, column is the class k (=1,2,...,K)	-- unsupervised -- Algo output: $Z - m \times K$ matrix (K <= n, and K=2 or 3 for easy plot)	$y - m \times 1$ vector; ( $y_i = 1$ anomalous; $y_i = 0$ normal)
Parameters	$\theta - (n+1) \times 1$ vector, $\theta_0$ is for bias feature	$\theta - (n+1) \times K$ matrix, column correspond to K output columns of y, row is parameter	Mapping function between layer L and layer L+1: $\Theta^{(L)} - S_{(L+1)} \times S_L + 1$ matrix	$\Theta - (m+1) \times K$ matrix, column correspond to K output columns of y, row is parameter	Null	$U_{reduced} - m \times n$ matrix	Null
Hypothesis	$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ $x_0 = 1$ (bias feature); $x_j$ can be $x^2, x^3, \dots$ or $x_0 x_1$ - a polynomial or nonlinear problem $h_{\theta}(x^{(i)}) = \tilde{\theta}^T x^{(i)}$ , $h_{\theta}(X) = X \cdot \theta$	$h_{\theta}^{(k)}(x) = P(y = k x; \theta)$ k = 1,2,3 ..., K $h_{\theta}(x) = g(\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n)$ $g(z) = \frac{1}{1+e^{-z}}$ - sigmoid/logistic function $h_{\theta}(x^{(i)}) = \frac{1}{1+e^{-\tilde{\theta}^T x^{(i)}}}$ , $h_{\theta}(X) = \frac{1}{1+e^{-X \cdot \theta}}$	 $Z^{(l+1)} = A^{(l)} \cdot \Theta^{(l)T}$ , $A^{(l+1)} = g(Z^{(l+1)})$ $A^{(l+1)}$ then add $m \times 1$ of "1" as bias units	Compute the kernels from given X: $f^{(i)} = x^{(i)}$ ; $f_m^{(i)} = \text{similarity}(x^{(i)}, x^{(m)}) = \exp\left(-\frac{\ x^{(i)} - x^{(m)}\ ^2}{2\sigma^2}\right)$ ; $f_0^{(i)} = 1$ (bias term) $h_{\theta}(f^{(i)}) = 1$ , if $\tilde{\theta}^T f^{(i)} \geq 0$ ; else $h_{\theta}(f^{(i)}) = 0$ Note: perform feature scaling before using the Gaussian Kernel. Other type of kernels: polynomial, string kernel, chi-square, histogram intersection, ...	1. Random initialization: <b>K cluster centroids</b> 2. Cluster assignment $c^{(i)} = \text{index}(1 \sim K)$ closest centroid to $x^{(i)}$ 3. Move centroids $\mu^{(k)} = \text{average of points assigned to cluster k}$ 4. compute cost function	1. Data processing: <b>mean normalization</b> 2. Compute "Covariance matrix": $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T - n \times n$ matrix 3. compute the eigenvectors of $\Sigma$ : [U, S, V] = svd(Sigma). Octave Code svd (singular value decomposition); Sigma -- $\Sigma$ 4. take the first K columns of $U_{n \times n}$ : $Z^{(i)} = U_{reduced}^T \cdot x^{(i)}$	Gaussian (Normal) distribution: $\text{Mean } \mu$ , $\text{variance } \sigma^2$ , $x \sim N(\mu, \sigma^2)$ $P(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$ Parameter Estimation: $\mu_1 = \frac{1}{m} \sum_{i=1}^m x_i^{(i)}$ , $\sigma_1^2 = \frac{1}{m} \sum_{i=1}^m (x_i^{(i)} - \mu_1)^2$ Multi-variant Gaussian distribution: $P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{1/2}  \Sigma ^{1/2}} \cdot \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$
Optimization Objective (cost function)	Mean Squared Error: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$ -- cost function of linear regression is always convex (converge) -- for sufficient small $\alpha$ , $J(\theta)$ should decrease on every iteration (plot $J(\theta)$ v.s. integration for debug)	Logistic cost function: $J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)}))]$	Feedforward propagation: $J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K [y_k^{(i)} \log h_{\theta}(x^{(i)})_k + (1 - y_k^{(i)}) \log (1 - h_{\theta}(x^{(i)})_k)]$ Sum over all K classes and all m training sets.	Cost function: $J(\theta) = C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\tilde{\theta}^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_2(\tilde{\theta}^T f^{(i)})] + \frac{1}{2} \sum_{i=1}^m \tilde{\theta}^2$ -- $C=1/\lambda$ , controls the penalty for misclassified training examples Using kernels with logistic regression can be very slow	Cost function: $J(C, \mu) = \frac{1}{m} \sum_{i=1}^m \ x^{(i)} - \mu_{c(i)}\ ^2$ -- repeat the entire steps (random initialization 100 times) to get rid of the local minimum; pick clustering that gave lowest cost function	1. Reconstruct with approximation $x_{approx}^{(i)} = U_{reduced} \cdot z^{(i)}$ 2. Choose the smallest K that can retain p (i.e. 99%) of variance: $\frac{\frac{1}{m} \sum_{i=1}^m \ x^{(i)} - \mu\ ^2}{\frac{1}{m} \sum_{i=1}^m \ x^{(i)}\ ^2} \leq 1 - p$	Given new example x, compute p(x): $p(x) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} \cdot \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$ Then predict: $y = \begin{cases} 1 & \text{if } p(x) < \epsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (normal)} \end{cases}$ Use CV to choose parameter $\epsilon$ -- data is very skew, we need to use precision/recall, F1-score.
Gradient Descent	$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$ $\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$ $\alpha$ -- learning rate	$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$	Backpropagation: $g'(z) = \frac{d}{dz} g(z) = g(z) \cdot (1 - g(z))$ $\delta^{(l)} = (\delta^{(l+1)} \cdot \Theta^{(l)}(:, 2:end)) \odot g'(z^{(l)})$ $\frac{\partial J(\theta)}{\partial \Theta^{(l)}} = \frac{1}{m} (\delta^{(l+1)})^T \cdot A^{(l)}$	-- SVM model computes the gradient for you	Null	Null	Null
Vectorization (regularized)	$J(\theta) = \frac{1}{2m} (X \cdot \theta - y)^T (X \cdot \theta - y) + \frac{\lambda}{2} \sum_{j=1}^n [\theta_{row1=0}^2]$ * here is sum over all matrix element squares $\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{m} X^T (X \cdot \theta - y) + \frac{\lambda}{m} \theta_{row1=0}$	$J(\theta) = -\frac{1}{m} \sum_{\text{along row}} [y \odot \log\left(\frac{1}{1+e^{-X \cdot \theta}}\right) + (1 - y) \odot \log\left(1 - \frac{1}{1+e^{-X \cdot \theta}}\right)] + \frac{\lambda}{2m} \sum_j [\theta_{row1=0}^2]$ , $1 \times K$ vector, $\odot$ is the element-wise product of two matrix $\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{m} X^T \left(\frac{1}{1+e^{-X \cdot \theta}} - y\right) + \frac{\lambda}{m} \theta_{row1=0}$	$J(\theta) = -\frac{1}{m} \sum [y \odot \log(h(\theta)) + (1 - y) \odot \log(1 - h(\theta))] + \frac{\lambda}{2m} \sum_{\text{layer}} \sum [\theta_{row1=0}^2]$ Regression part: sum over all matrix elements; regularization part: sum over all matrix elements and all layer $\Theta$ matrix $\frac{\partial J(\theta)}{\partial \Theta^{(l)}} = \frac{1}{m} (\delta^{(l+1)})^T \cdot A^{(l)} + \frac{\lambda}{m} \Theta_{col,1=0}$ "unroll" each $\Theta$ and append together	$J(\theta) = C \sum_{\text{along row}} [y \odot \text{cost}_1(\tilde{\theta}^T f^{(i)}) + (1 - y) \odot \text{cost}_2(\tilde{\theta}^T f^{(i)})] + \frac{1}{2} \sum_{i=1}^m [\theta_{row1=0}^2]$ -- large C: lower bias, high variance -- large $\sigma^2$ : features $f^{(i)}$ vary more smoothly, high bias and low variance	Null	$Z = X \cdot U_{reduced}$ Compute the % of variance from diagonal matrix S: $S = \begin{bmatrix} S_{11} & 0 & 0 \\ 0 & \dots & 0 \\ 0 & 0 & S_{nn} \end{bmatrix}$ , $\frac{\sum_{j=1}^K S_{jj}}{\sum_{j=1}^n S_{jj}} \geq p$	Transform non-Gaussian features to make them look more Gaussian: E.g.1: $x_j = \log(x_j + c)$ E.g.2: $x_j = \sqrt{x_j + c}$ Multi-variant Gaussian Distribution parameter: $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$ ; $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$
Decision boundary	Null	$\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = 0$ or $\tilde{\theta}^T x^{(i)} = 0$ ( $\tilde{\theta}^T x^{(i)} > 0, y = 1$ ; else $\tilde{\theta}^T x^{(i)} < 0, y = 0$ )	Null	$\tilde{\theta}^T f^{(i)} = 0$ , but with the choice of $\theta$ that: $\tilde{\theta}^T f^{(i)} \geq 1$ if $y^{(i)} = 1$ ; $\tilde{\theta}^T f^{(i)} \leq -1$ if $y^{(i)} = 0$ effectively more margin	Null	Null	Null
Notes /Limitations	-- Pre-process data by <b>mean normalization</b> of each feature: $x_j = \frac{x_j - \mu_j}{s_j}$ where $\mu_j$ is mean and $s_j$ is range/ standard deviation (save computation time) -- analytical solution of $\theta$ when dimension (of training set is small): $\theta = (X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \dots \end{bmatrix})^{-1} X^T y$	-- $h_{\theta}(x)$ compute the probability of each class; for the new test set, output is the class k that maximize the $h_{\theta}(x)$ -- binary classification/ multi-class classification	-- randomly initialize $\Theta_{ij}^{(l)}$ to a random value in [-e, +e] (symmetry breaking to avoid redundant hidden units) -- we can choose $\epsilon_{\text{initial}} = \frac{\sqrt{6}}{\sqrt{S_L + S_{L+1}}}$ -- usually more hidden layer is better, choose same # of units in every hidden layer	-- Idea of kernels is to set each training example as a landmark $l^{(i)} = x^{(i)}$ , then transfer each $x^{(i)} \rightarrow f^{(i)}$ , wherein each element $f_m^{(i)} = \text{similarity}(x^{(i)}, x^{(m)})$ -- if n is small and m is intermediate (n=1~1000, m=10~1000), we can use SVM with kernel; otherwise, use logistic regression or SVM without kernel -- for a linear classification problem, use logistic regression or SVM without kernel	-- suitable for K not too large (2~10) -- choose cluster number K: (1) elbow method (i.e. J v.s. K plot); (2) based on later/downstream purpose (i.e. T-shirt size)	-- Reduce from n-dimension to k-dimension: find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error -- This mapping optimization should be defined by running PCA only on training set, and then applied to CV or test set. -- <b>misuse</b> includes: (1) use PCA to reduce features to prevent overfitting; (2) use PCA before even try with original data set	-- Differ from supervised learning, anomaly detection has very small # of positive examples, hard for any algo to learn the positive examples, and new anomalies may look nothing like trained examples. -- Choose feature that might take on unusually large/small values (exaggerate the anomaly signal) -- If features form clear correlation, use multivariate GD in case of miss some anomaly ( $\Sigma$ has off-diagonal elements)
Applications	# linear regression problem # recommender system	# Multi-class Classification: One-vs-All	# Binary operations (i.e. XOR) # Multi-class Classification: One-vs-All	# non-linear Multi-class Classification # One-vs-All (e.g. spam-email classification, weather prediction)	# Market segmentation; # social network analysis; # organize computing clusters; # Astronomical data analysis	# image compression # data visualization # reduce dimension to speed up learning algorithm	# Fraud detection # Manufacturing outliers # Monitor

## Build a Machine Learning System

ML System	Algorithm Computation				ML Diagnostic													
	Optimization Algorithm	Regularization	Gradient Checking	Evaluate a Hypothesis	Bias/Variance & Learning Curve	Error Analysis	Ceiling Analysis											
Explanations	<p><i>Basic:</i></p> <p>-- Gradient descent</p> <p><i>Advanced:</i></p> <p>-- Conjugate gradient; BFGS; L-BFGS</p> <p><i>Note:</i></p> <p>-- advanced algos don't need to pick the learning rate, often faster than gradient descent</p>	<p><i>Motivation:</i></p> <p>Address overfitting issue by reduce the magnitudes of the parameters <math>\theta_j</math> while keep the # of features unchanged</p>	<p><i>Motivation:</i></p> <p>Numerically estimate the gradient in order to validate the gradient decent terms.</p> <p>Disable it before running the advanced optimization algorithm</p>	<p><i>Evaluate a Hypothesis:</i></p> <p>1. shuffle data set and divide into 3 groups: 60% training; 20% cross-validation (CV); 20% test</p> <p>2. <u>Optimizing the parameter based on training error -&gt; determine the degree of model based on CV error -&gt; estimate the generalization error with test set</u></p>	<p><u>Under fitting (high bias):</u></p> <p>Both <math>J_{train}</math> and <math>J_{CV}</math> are high; <math>\lambda</math> could be too large; Plot learning curve (error v.s. training set size <math>m</math>). <math>J_{train}</math> and <math>J_{CV}</math> converge but both very high</p> <p><u>Over fitting (high variance):</u></p> <p><math>J_{train}</math> is low and <math>J_{CV}</math> is high; <math>\lambda</math> could be too small; Plot learning curve, <math>J_{train}</math> and <math>J_{CV}</math> converge with a gap; increase training set size can help</p>	<p><i>Motivation: Rare class/skew class</i></p> <p>Training example belongs to one class out-numbers heavily the other classes.</p> <table><tr><td colspan="2"></td><td colspan="2">Actual Class</td></tr><tr><td rowspan="2">Predict class</td><td>1</td><td>True Positive</td><td>False Positive</td></tr><tr><td>0</td><td>False Negative</td><td>True Negative</td></tr></table>			Actual Class		Predict class	1	True Positive	False Positive	0	False Negative	True Negative	<p><i>Motivation: Prioritize the high-error modules in a ML pipeline</i></p> <p>Machine learning pipeline:</p>
			Actual Class															
Predict class	1	True Positive	False Positive															
	0	False Negative	True Negative															
Applications	<p><i>Octave Code</i></p> <p>Step1: write a function</p> <p><math>J, \text{grad} = \text{costFunction}(\text{theta}, X, y)</math></p> <p>Step2: compute the optimal theta</p> <p><i>Options = optimset('GradObj','on','MaxIter','100');</i></p> <p><i>initialTheta = zeros(2,1);</i></p> <p><i>[optTheta, functionVal, exitFlag] = fminunc(@costFunction, initialTheta, options);</i></p>	<p><i>Cost function:</i></p> $J(\theta) = J(\theta) + \frac{\lambda}{2m} \sum_{j=1}^n \sum_{k=1}^K \theta_{j,k}^2$ <p>* Remember <math>\theta</math> is <math>(n+1) \times k</math>, the first row of <math>\theta_j</math> are excluded</p> <p><i>Gradient:</i></p> $\frac{\partial J(\theta)}{\partial \theta} = \frac{\partial J(\theta)}{\partial \theta} + \frac{\lambda}{m} \theta_{\text{first\_row}=0}$	<p><i>Concept:</i></p> $\frac{\partial J(\theta)}{\partial \theta_j} \approx \frac{J(\theta_j + \epsilon, \theta_{-j}) - J(\theta_j - \epsilon, \theta_{-j})}{2\epsilon}$ <p><i>Octave Code</i></p> <pre>For i=1:n,     thetaPlus = theta; thetaMinus = theta;     thetaPlus += epsilon; thetaMinus = epsilon;     gradApprox(i) = (J(thetaPlus) - J(thetaMinus))/2/epsilon; end;</pre> <p>check that gradApprox <math>\approx</math> DVec from backpropagation (neuralNetwork)</p>	<p><i>Test set error:</i></p> <p>Regression problem – mean square error (same as cost function)</p> <p>Classification problem: <math>\text{err}(h_\theta(x), y) = 1</math> if <math>h_\theta(x) \geq 0.5</math> and <math>y=0</math> or if <math>h_\theta(x) &lt; 0.5</math> and <math>y=1</math>; otherwise, <math>\text{err}(h_\theta(x), y) = 0</math></p> <p><i>Note:</i></p> <p>When optimizing the parameters <math>\theta</math> by training set, use regularization.</p> <p>When plotting training/CV error, do not include regularization.</p>	<p>-- "small" neural network (fewer hidden layers/units) is more prone to underfitting (high bias).</p> <p>-- high bias and high variant can both happen, as regional over/under fit</p> <p><b>Build ML system:</b></p> <ol style="list-style-type: none"><li>1. start with a simple algorithm</li><li>2. plot learning curve</li><li>3. Error analysis on CV set (determine if high bias or variance)</li></ol>	<p><i>Equations:</i></p> $\text{precision} = \frac{tp}{tp+fp}; \text{recall} = \frac{tp}{tp+fn}$ $F1 \text{ score} = \frac{2PR}{P+R}$ <p><b>Threshold:</b></p> <p>For logic regression: <math>0 \leq h_\theta(x) \leq 1</math></p> <p>If <math>h_\theta(x) \geq \text{threshold}</math>, predict 1.</p> <p><i>Note:</i></p> <p>-- increase threshold can improve the precision but trade off the recall</p>	<p>E.g. manually feed perfect "text detection" system accuracy can be improved to 89%. Improve "text detection" can bring 17% up</p> <p>Improve "Character Segm" only bring 1% up</p> <table><tr><th>Component</th><th>Accuracy</th></tr><tr><td>Overall system</td><td>72%</td></tr><tr><td>→ Text detection</td><td>89%</td></tr><tr><td>Character segmentation</td><td>90%</td></tr><tr><td>Character recognition</td><td>100%</td></tr></table>	Component	Accuracy	Overall system	72%	→ Text detection	89%	Character segmentation	90%	Character recognition	100%	
Component	Accuracy																	
Overall system	72%																	
→ Text detection	89%																	
Character segmentation	90%																	
Character recognition	100%																	

## Special Applications

	Concept/Model	Computation	Notes	Applications
Recommender Systems	<p>Problem Statement: Given output, while input and parameters are not defined</p> <p>Model: <b>Linear regression model</b></p> <p>Output:</p> <p><math>Y = m \times u</math> matrix, <math>y_i</math> – user <math>j</math>'s rating on movie <math>i</math> (row is movie, column is user)</p> <p><math>R = m \times u</math> matrix, <math>r_{ij} = 1</math> – user <math>j</math> has rated movie <math>i</math> (otherwise <math>r_{ij} = 0</math>)</p> <p>Parameters:</p> <p><math>X = m \times (n+1)</math> matrix, <math>x_{ik}</math> – movie <math>i</math> has feature <math>k</math> (i.e. romance, action, etc.)</p> <p><math>\theta = (n+1) \times u</math> matrix, <math>\theta_{ij}</math> – user <math>j</math>'s parameter for feature <math>k</math></p>	<p>Algorithm: <b>Collaborative Filtering Algorithm</b></p> <p>Cost function:</p> $\min_{\theta, x} \frac{1}{2} \sum_{(i,j) \in R} \left( (\theta^{(j)})^T \cdot x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n \left( \theta_k^{(j)} \right)^2 + \frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n \left( x_k^{(i)} \right)^2$ <p>Vectorization: <math>J(\theta, x) = \frac{1}{2} \sum (X\theta\Theta R - Y\Theta R)^2 + \frac{\lambda}{2} \sum \sum \theta^2 + \frac{\lambda}{2} \sum \sum x^2</math></p> <p>Gradient: <math>\frac{\partial J(\theta, x)}{\partial \theta} = X^T (X\theta\Theta R - Y\Theta R) + \lambda \theta_{\text{column}=0}</math></p> $\frac{\partial J(\theta, x)}{\partial x} = X^T (X\theta\Theta R - Y\Theta R) + \lambda x_{\text{row}=0}$ <p>Find related movies: <math>\ x(i) - x(j)\  &lt; \epsilon</math></p>	<p>-- In this problem, both parameters <math>x</math> and <math>\theta</math> are unknown. Putting both <math>\theta</math> and <math>x</math> optimization objective together and combine into a single cost function, instead of optimizing <math>\theta</math> and <math>x</math> back-forth in serial (<math>\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \dots</math>)</p> <p>-- For users who have not rated any movies, use the mean normalization of the training examples (existing users)</p>	# recommend movies/ads for users
Large Scale Machine Learning	<p><i>Before deploying expensive ML algorithm:</i></p> <ul style="list-style-type: none"> <li>-- Use relatively small data set -&gt; plot learning curve -&gt; confirm it's high variant</li> <li>-- Understand what's cost to simply gather more data</li> </ul> <p><b>Map Reduce and data parallelism:</b></p> <ul style="list-style-type: none"> <li>-- Reduce the training set into several subset</li> <li>-- Compute the training set separately (i.e. in different computers or CPU cores) and combine them up</li> </ul>	<p><b>"Batch" gradient descent:</b> <math>\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}</math></p> <p><b>Stochastic gradient descent:</b> randomly shuffle data set, then <math>\frac{\partial J(\theta)}{\partial \theta_j} = (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}</math></p> <p><b>"Mini-batch" gradient descent:</b> <math>\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{b} \sum_{i=1}^b (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}</math></p>	<p>-- For stochastic gradient, the learning rate <math>\alpha</math> is typically held constant, we can slowly decrease <math>\alpha</math> over time if we want <math>\theta</math> to converge (E.g. <math>\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}</math>)</p>	# online learning (continuous flood of data, compute gradient only on new data, can adapt to user behavior changes)