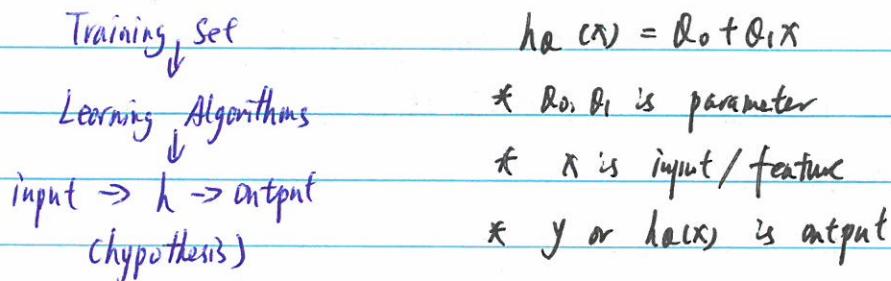


DS01 Machine Learning (Stanford)

Week 1

Introduction and Univariate Linear Regression

- Algorithms:
 - Supervised learning / Unsupervised learning / Reinforcement learning / Recommender system
 - Supervised learning: "right answer" given
 - Regression: predict continuous valued output
 - Classification: discrete valued output (0 or 1)
- Unsupervised learning:
 - Kmean, cluster
- Linear regression w one variable
- Model representation:



- Cost function: $J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$ (squared error function)
 Goal: $\min_{(\theta_0, \theta_1)} J(\theta_0, \theta_1)$ in training set $h_{\theta}(x_i) = \theta_0 + \theta_1 x^i$
 - Gradient descent algorithm
 - repeat until converge: $\left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \right\}$
 - simultaneously update θ_0 and θ_1
- $$\text{temp } 0 = \theta_0 - \alpha \frac{d}{d\theta_0} J(\theta_0, \theta_1) = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)$$
- $$\text{temp } 1 = \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_0, \theta_1) = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) x_i$$
- $$\theta_0 = \text{temp } 0$$
- $$\theta_1 = \text{temp } 1$$
- update θ_0 and θ_1 simultaneously
- * α too small \Rightarrow too slow to converge, α too large \Rightarrow fail converge

- * Gradient descent can converge to local min, even with fixed α
As we approach to local min, gradient descent will automatically take smaller step, so no need to decrease α over the time
- * "Batch": each step of gradient descent use all training example

Week 2

Multivariate Linear Regression

- Linear regression with multiple variables
- Hypothesis: $h_{\theta}(x) = \theta^T x$, parameter $\theta = [\theta_0 \ \dots \ \theta_n]$, input $x = [x_0 \ \dots \ x_n]$
- Cost function: $J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$
- Gradient descent: Repeat $\{\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)\}$
Simultaneously update θ_j : $\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) x_i$
- Feature Scaling:
 - | make sure feature are on similar scale
 - | give every feature in $-1 \leq x_i \leq 1$
- Mean normalization: $x_i = (x_i - \bar{x}_i) / s_i$ where $s_i = \text{range}(\text{max-min})$ or std
- Learning rate
 - $J(\theta)$ should decrease after every iteration
 - Declare convergence if $J(\theta)$ decrease by $< 10^{-3}$ in 1 iteration
- for smaller α , $J(\theta)$ should decrease but can be slow to converge
for larger α , $J(\theta)$ may not decrease and may not converge
- Normal equation: solve θ analytically \leftarrow work well if small x .
 - $\frac{\partial}{\partial \theta_j} J(\theta) = 0$ solve $\theta_0 - \theta_n$ where $J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$
 - $\theta = (X^T X)^{-1} X^T y$
 - $(X^T X)$ is non-invertible if
 - | Redundant feature (linear dependent)
 - | too many features, delete some features or use regularization

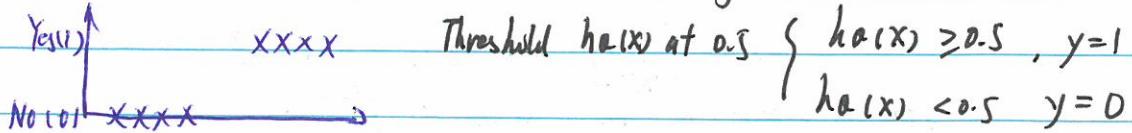
Week 3

Logistic Regression

- Classification

- $y \in \{0, 1\}$ 0: Negative class (benign tumor)

1: Positive class (malignant tumor)



- Classification: $y=0$ or 1

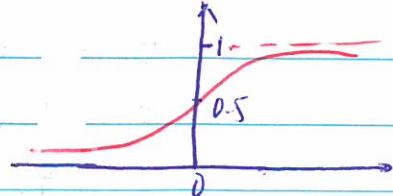
Logistic regression: $0 \leq h_0(x) \leq 1$

• Hypothesis Representation

- Logistic Regression Model

$$0 \leq h_0(x) \leq 1$$

$$\begin{aligned} h_0(x) &= g(\theta^T x) \\ g(z) &= \frac{1}{1+e^{-z}} \end{aligned} \quad \Rightarrow \quad h_0(x) = \frac{1}{1+e^{-\theta^T x}}$$



* Sigmoid function / Logistic function

- Interpretation of Hypothesis Output

* $h_0(x)$ = estimated probability that $y=1$ on input x

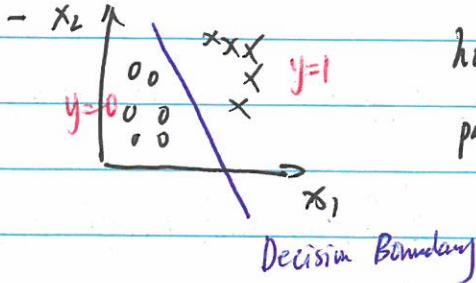
* $h_0(x) = P(y=1 | x; \theta)$ probability that $y=1$, given x , parameterized by θ

$$P(y=0 | x; \theta) = 1 - P(y=1 | x; \theta)$$

• Decision Boundary

- $h_0(x) = g(\theta^T x) \Rightarrow y=1 \text{ if } h_0(x) \geq 0.5 \Rightarrow \theta^T x \geq 0$

$$g(z) = \frac{1}{1+e^{-z}} \quad y=0 \text{ if } h_0(x) < 0.5 \Rightarrow \theta^T x < 0$$



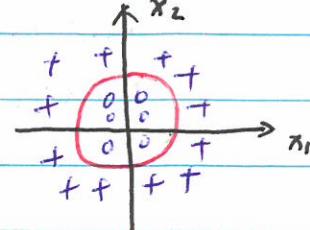
$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

predict $y=1$ if $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

- Nonlinear Decision Boundary

$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

predict $y=1$ if $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 \geq 0$



- Cost Function

- Training set: $\{(x^1, y^1), \dots, (x^m, y^m)\}$

$$x \in \begin{bmatrix} x^1 \\ \vdots \\ x^m \end{bmatrix}, y_i = 1, y \in \{0, 1\}$$

How to choose θ in $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$

- $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^i), y^i)$ where $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$

* non-convex

$$\text{Cost}(h_\theta(x), y)$$

$$\text{Cost} \uparrow$$

$$y=1$$

$$y=0$$

- $\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$

* Cost $\rightarrow \infty$ if $y=1, h_\theta(x)=1$, as $h_\theta(x) \rightarrow 0$. Cost $\rightarrow \infty$

* Captures intuition that if $h_\theta(x)=0$, (predict $P(y=1|x; \theta) \approx 0$), but

$y=1$, we'll penalize learning algorithm by a very large cost

- Simplified Cost Function and Gradient Descent

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

- $\text{Cost}(h_\theta(x^i), y^i) = \begin{cases} -\log(h_\theta(x^i)) & \text{if } y^i=1 \\ -\log(1-h_\theta(x^i)) & \text{if } y^i=0 \end{cases}$

- $\text{Cost}(h_\theta(x^i), y^i) = -y^i \log(h_\theta(x^i)) - (1-y^i) \log(1-h_\theta(x^i)) \quad y^i=0 \text{ or } 1$

if $y^i=1$: $\text{Cost}(h_\theta(x^i), y^i) = -\log(h_\theta(x^i))$

if $y^i=0$: $\text{Cost}(h_\theta(x^i), y^i) = -\log(1-h_\theta(x^i))$

$\Rightarrow \left\{ \begin{array}{l} \text{Cost function } J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^i \log h_\theta(x^i) + (1-y^i) \log(1-h_\theta(x^i)) \right] \\ \text{To find parameter } \theta: \min_{\theta} J(\theta) \end{array} \right.$

To predict new x : output $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$ $P(y=1|x; \theta)$

- Gradient Descent:

Want $\min_{\theta} J(\theta)$ where $J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^i \log h_\theta(x^i) + (1-y^i) \log(1-h_\theta(x^i)) \right]$

Repeat { $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ } Simultaneously update all θ_j

where $\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i) x_j^i$

$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$

- Advanced Optimization

- Cost function $J(\alpha)$, want $\min_{\alpha} J(\alpha)$

Given α , we compute $J(\alpha), \frac{\partial}{\partial \alpha_j} J(\alpha)$

Gradient descent = Repeat $\{ \alpha_j := \alpha_j - \alpha \frac{\partial}{\partial \alpha_j} J(\alpha) \}$

- Optimization algorithms.

* Conjugate gradient / BFGS / L-BFGS

* Advantage: No need to pick α / much faster

* options = optimset('GradObj', 'on', 'MaxIter', '100')

- Multi-class classification: One-vs-all

- One-vs-all = one-vs-rest

$$h_\alpha^i(x) = P(y=i|x; \alpha)$$

* Train a logistic regression $h_\alpha^i(x)$ for each class i to predict the probability that $y=i$

- The Problem of Overfitting

- Reduce # of features: manually / Model selection algorithm

- Regularization: keep all features, but reduce α_j

- Regularization Cost Function

- Small values of α_j : simpler hypothesis / less prone to overfitting

- $J(\alpha) = \frac{1}{m} \sum_i^m (h_\alpha(x^i) - y^i)^2 + \lambda \sum_j^N \alpha_j^2$ regularization parameter

* $\lambda \uparrow$, algorithm works fine / fails to eliminate overfitting / results in underfitting / gradient descent fail to converge

- Regularized Linear Regression

- Repeat $\{ \alpha_0 := \alpha_0 - \alpha \frac{1}{m} \sum_i^m (h_\alpha(x^i) - y^i) x_0^i \}$

$$\alpha_j := \alpha_j (1 - \alpha \frac{1}{m}) - \alpha \frac{1}{m} \sum_i^m (h_\alpha(x^i) - y^i) x_j^i$$

- Normal Equation: $\alpha = (X^T X + \lambda [1, \dots, 1]^T)^{-1} X^T y$
 $(m+1) \times (n+1)$

- Regularized Logistic Regression

$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y_i \log h_\theta(x_i) + (1-y_i) \log (1-h_\theta(x_i)) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

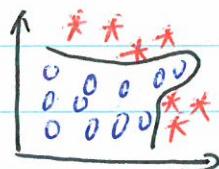
$$\text{Repeat: } \begin{cases} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) x_i \\ \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) x_i^j - \frac{\lambda}{m} \theta_j \right] \end{cases}$$

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) x_i^j - \frac{\lambda}{m} \theta_j \right]$$

Week 4

Neural Networks I

• Nonlinear Hypotheses:

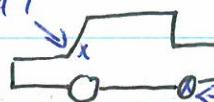


$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 x_2 + \dots + \theta_n x_1 x_n)$$

$O(n^2)$

- Computer Vision: Car detection (car vs not a car)

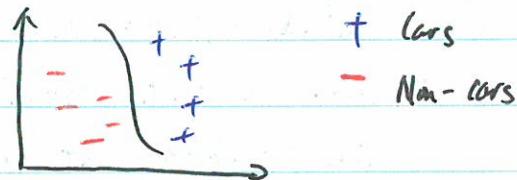
pixel 1



50x50 pixel image \rightarrow 25w pixels

$$\mathbf{x} = \begin{bmatrix} \text{pixel 1} \\ \vdots \\ \text{pixel 25w} \end{bmatrix}$$

\Rightarrow Quadratic features ($x_i x_j$): 3 million features



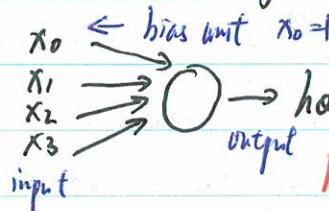
+ cars

- Non-cars

- Neurons and the brain \Leftrightarrow Neural Network

- Model Representation

- Neuron model: Logistic unit

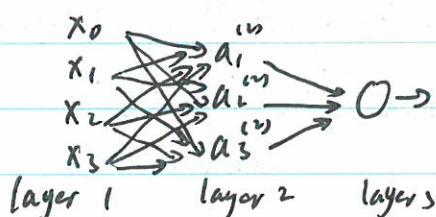


$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$\mathbf{x}_2 = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} \quad (\text{weight})$$

$$\text{Sigmoid (logistic) activation function } g(z) = \frac{1}{1 + e^{-z}}$$

- Neural Network



$a_i^{(j)}$ = activation of unit i in layer j

$\boldsymbol{\theta}^{(j)}$ = matrix of weight from layer j to $j+1$

$$a_1^{(1)} = g(Q_{10}^{(1)}x_0 + Q_{11}^{(1)}x_1 + Q_{12}^{(1)}x_2 + Q_{13}^{(1)}x_3)$$

$$a_2^{(1)} = g(Q_{20}^{(1)}x_0 + Q_{21}^{(1)}x_1 + Q_{22}^{(1)}x_2 + Q_{23}^{(1)}x_3)$$

$$h(x) = a_1^{(1)} = g(Q_{10}^{(1)}a_0^{(0)} + Q_{11}^{(1)}a_1^{(0)} + Q_{12}^{(1)}a_2^{(0)} + Q_{13}^{(1)}a_3^{(0)})$$

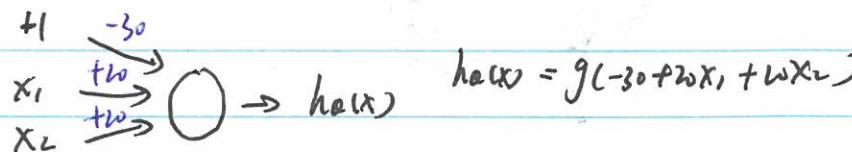
* if network has S_j unit in layer j , S_{j+1} units in layer $j+1$, then $Q^{(j)}$ will be of dimension $S_{j+1} \times (S_j + 1)$

- Forward propagation : Vectorized implementation

$$z^{(1)} = Q^{(0)}a^{(0)} \Rightarrow a^{(1)} = g(z^{(1)}) \Rightarrow z^{(2)} = Q^{(1)}a^{(1)} \Rightarrow h(x) = a^{(2)} = g(z^{(2)})$$

- Example,

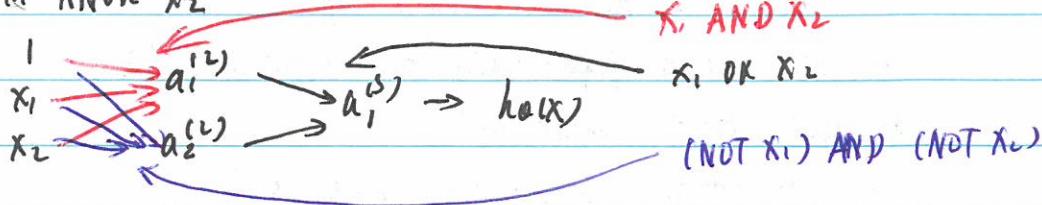
- AND : $x_1, x_2 \in \{0, 1\}$, $y = x_1 \text{ AND } x_2$



- OR : $x_1, x_2 \in \{0, 1\}$, $y = x_1 \text{ OR } x_2$

- NOT x_1

- $x_1 \text{ XNOR } x_2$



- Multi-class classification

- One-vs-all : $h(x) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ $h(x) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ $h(x) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

when pedestrian when car when motorcycle

Train set : $(x^{(1)}, y^{(1)})$, $(x^{(2)}, y^{(2)})$... $(x^{(m)}, y^{(m)})$

where $y^{(i)}$ one of $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$
pedestrian car motorcycle truck

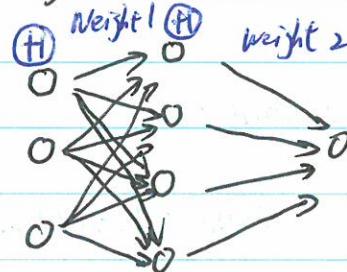
- How to build your own Neural Network from scratch in Python

- What is a Neural Network

* Neural Network consists of

- o Input layer, x
- o Arbitrary amount of hidden layer
- o Output layer, \hat{y}
- o A set of weights and biases between each layer, w and b { Similar to θ in the course
- o Activation function for each hidden layer σ (Sigmoid activation function)

* 2-layers neural network



Input Layer Hidden Layer Output layer

- Training the Neural Network

$$\hat{y} = \sigma(w_2 \sigma(w_1 x + b_1) + b_2)$$

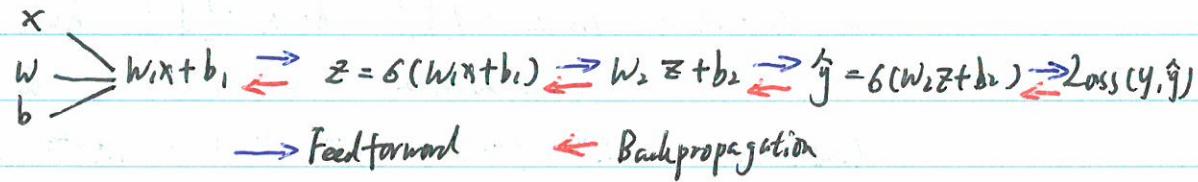
* w and b are the only variables affect \hat{y}

* we need to fine-tuning the weight and biases

* each iteration consists of

o feedforward : calculating predicted \hat{y}

o backpropagation : updating weight and biases



- FeedForward

$$\hat{y} = \sigma(w_2 \sigma(w_1 x + b_1) + b_2)$$

- Loss function

$$\text{Sum-of-square error} = \frac{1}{N} \sum_i (y - \hat{y})^2$$

* or $-\frac{1}{m} \sum_{i=1}^m [y_i \log a_i^L + (1-y_i) \log (1-a_i^L)]$

- * Goal : find best set of weights and biases to min loss function
 - Backpropagation

* Use Loss function to measure the error of prediction, and then propagate the error back and update weight and biases

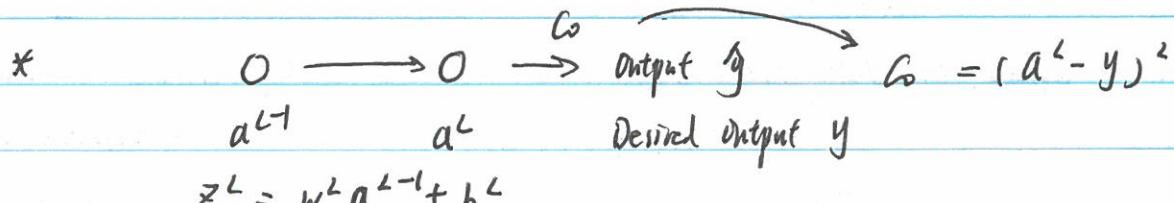
- o Derivative of loss function to weight and biases
- o and use gradient descent to update weight and biases

* Use Chain rule to get Derivative

$$\text{Loss}(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$$

$$\frac{d \text{Loss}(y, \hat{y})}{d w} = \frac{d \text{Loss}(y, \hat{y})}{d \hat{y}} * \frac{d \hat{y}}{d z} * \frac{d z}{d w} \text{ where } z = w \cdot x + b$$

$= 2(y - \hat{y}) * \text{derivative of sigmoid function}$

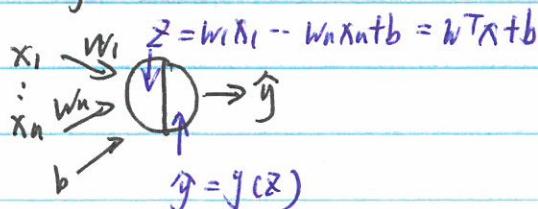


$$\hat{y} = a^L = g(z^L)$$

- Deep Neural Network from scratch in Python

- a_i^L L is L th layer, i is i th example

- Single neuron



* A neuron calculate $z = w^T x + b$

followed by an activation function

* g is activation function (Sigmoid, tanh, ReLU)

- Work flow :

1. normalized all feature : to center and standardize dataset $\frac{x - \bar{x}}{\sigma}$
 define the inputs

2. Initialize parameters and define hyperparameters:

- * # of iteration

- * # of layers L in the neural network

- * size of hidden layer

- * Learning rate α

3. Loop for num_iteration:

- * Forward propagation (Feed forward) to calc current loss

- * Compute cost function

- * Backward propagation (Back propagation) to calc current gradient

- * Update parameters

4. Use trained parameter to predict labels

- Activation functions

Sigmoid: * def sigmoid(z):

$$S = 1 / (1 + \exp(-z))$$

return S

def sigmoid_backward(dA, z)

$$S = \text{sigmoid}(z)$$

$$dS = S * (1 - S)$$

return dS * dA

ReLU: * def relu(z):

$$R = \text{np. maximum}(0, z)$$

return R

def relu_backward(dA, z)

$$dz = \text{np. array}(dA, \text{copy}=\text{True})$$

$$dz[z <= 0] = 0$$

return dz

- Forward propagation

- * $z^l = W^l A^{l-1} + b^l$

$$A^l = g(z^l)$$

- Loss function

- * $-\frac{1}{m} \sum_{i=1}^m [y_i \log a_i^L + (1-y_i) \log (1-a_i^L)]$

- Backward propagation

- * using chain rule

- Update parameter: $\text{parameter}_i = \text{parameter}_i - \text{learning rate} * \text{grads}$

- To avoid overfitting

* use regularization:

o L2 regularization: appropriately modify cost function

o drop out: randomly shuts down some neurons in each iteration

Week 5 Neural Network 2 and Practical Advice

- Neural Network Learning - Cost Function

- Initialization

$$\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) - (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}$$

L = total # of layers

s_l = # of units (not include bias unit) in layer l

* Binary classification: $y = 0$ or 1

* Multi-class classification (K class): $[1]_0$ car $[0]_1$ truck ---

- Cost function:

$$\text{Logistic regression: } J(\boldsymbol{\alpha}) = -\frac{1}{m} \left[\sum_{i=1}^m y_i \log h_{\boldsymbol{\alpha}}(\mathbf{x}^i) + (1-y_i) \log (1-h_{\boldsymbol{\alpha}}(\mathbf{x}^i)) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \beta_j^2$$

$$\text{Neural network: } J(\boldsymbol{\alpha}) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^i \log h_{\boldsymbol{\alpha}}(\mathbf{x}^i)_k + (1-y_k^i) \log (1-h_{\boldsymbol{\alpha}}(\mathbf{x}^i)_k)$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\alpha_{ji}^l)^2$$

- Neural Networks Learning - Backpropagation Algorithm

$$\text{- Goal: } J(\boldsymbol{\alpha}) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^i \log (h_{\boldsymbol{\alpha}}(\mathbf{x}^i)_k) + (1-y_k^i) \log (1-h_{\boldsymbol{\alpha}}(\mathbf{x}^i)_k) \right]$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\alpha_{ji}^l)^2$$

$$\min J(\boldsymbol{\alpha})$$

need to calculate: $J(\boldsymbol{\alpha})$ and $\frac{\partial}{\partial \alpha_{ij}^l} J(\boldsymbol{\alpha})$

- Forward propagation:

Given one training example (\mathbf{x}, \mathbf{y})

$$\mathbf{a}^1 = \mathbf{x}$$

Layer	1	2	3	4
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

$$\mathbf{a}^1 \quad \mathbf{a}^2 \quad \mathbf{a}^3 \quad \mathbf{a}^4$$

$$z^1 = \theta^0 a^0$$

$$a^1 = g(z^1) \quad \text{add } a_0^1 \quad (\text{usually } = 1)$$

$$z^2 = \theta^1 a^1$$

$$a^2 = g(z^2) \quad \text{add } a_0^2$$

$$z^3 = \theta^2 a^2$$

$$a^3 = g(z^3) \quad \text{add } a_0^3$$

$$z^4 = \theta^3 a^3$$

$$a^4 = h(\theta(x)) = g(z^4)$$

- Backpropagation Algorithm

* δ_j^l = "error" of node j in layer l

$$\text{Chain Rule} \downarrow \quad \delta_j^4 = a_j^4 - y_j = (h(\theta(x)))_j - y_j$$

$$\delta^3 = (\theta^3)^T \delta^4 \cdot g'(z^3) \leftarrow \text{if } g(\cdot) = \text{sigmoid}, g'(z^3) = a^3 \cdot (1-a^3)$$

$$\delta^2 = (\theta^2)^T \delta^3 \cdot g'(z^2) \leftarrow a^2 \cdot (1-a^2)$$

$$\frac{\partial}{\partial \theta_{ij}^l} J(\theta) = a_j^{l+1} \delta_i^{l+1}$$

* Training Set $\{(x^1, y^1), \dots, (x^m, y^m)\}$

$$\text{Set } \Delta_{ij}^l = 0 \quad \leftarrow \text{Used to update } \frac{\partial}{\partial \theta_{ij}^l} J(\theta)$$

for $i=1$ to m :

$$\text{set } a^1 = x^i$$

perform forward propagation to compute a^l for $l=2, 3, \dots, L$

using y^i , compute $\delta^L = a^L - y^i$

compute $\delta^{L-1}, \delta^{L-2}, \delta^{L-3}, \dots, \delta^2$

$$\Delta_{ij}^l := \Delta_{ij}^{l-1} + a_j^{l+1} \delta_i^{l+1} \leftarrow (\Delta^l := \Delta^l + \delta^{l+1} (a^l)^T)$$

$$D_{ij}^l := \frac{1}{m} \Delta_{ij}^l + \lambda \theta_{ij}^l \quad \text{if } j \neq 0$$

$$D_{ij}^l := \frac{1}{m} \Delta_{ij}^l \quad \text{if } j = 0$$

$$\frac{\partial}{\partial \theta_{ij}^l} J(\theta) = D_{ij}^l$$

* δ_j^l = "error" of cost for a_j^l (unit j in layer l)

Formally, $\delta_j^l = \frac{\partial}{\partial z_j^l} \text{cost}(i)$, where $\text{cost}(i) = y^i \log a^i + (1-y^i) \log (1-a^i)$

- Implementation Note:

* Advanced optimization

Layer	1	2	3	4
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
δ^2	δ^3	δ^4		
a^1	a^2	a^3	a^4	

function ZjVal, gradient J = cost Function (theta)

opt Theta = fminunc(@(cost Function, initialTheta, options)

$\Omega^1, \Omega^2, \Omega^3$ - matrices (Theta1, Theta2, Theta3) $\vec{R}^{10 \times 11}, \vec{R}^{10 \times 11}, \vec{R}^{1 \times 11}$

D^1, D^2, D^3 - matrices (D1, D2, D3) $\vec{R}^{10 \times 11}, \vec{R}^{10 \times 11}, \vec{R}^{1 \times 11}$

* From thetaVec, get $\Omega^1, \Omega^2, \Omega^3$

* Use forward/backpropagation to compute D^1, D^2, D^3 and $J(\Omega)$. Unroll

D^1, D^2, D^3 to get gradient Vec

- Gradient check

$$\frac{d}{d\Omega} J(\Omega) = J(\Omega + \epsilon) - J(\Omega - \epsilon) / 2\epsilon$$

$$\text{gradApprox} = (J(\theta + \epsilon \text{EPSILON}) - J(\theta - \epsilon \text{EPSILON})) / (2 * \epsilon \text{EPSILON})$$

$$\Omega = [\Omega_1 \dots \Omega_n]$$

$$\frac{d}{d\Omega_i} J(\Omega) \approx \frac{J(\Omega_i + \epsilon, \Omega_1 \dots \Omega_n) - J(\Omega_i - \epsilon, \Omega_1 \dots \Omega_n)}{2\epsilon}$$

$$\frac{d}{d\Omega_n} J(\Omega) \approx \frac{J(\Omega_1 \dots \Omega_{n-1}, \Omega_n + \epsilon) - J(\Omega_1 \dots \Omega_{n-1}, \Omega_n - \epsilon)}{2\epsilon}$$

check that gradApprox \approx DVec

- Neural Network Learning - Random Initialization

* Initial each Ω_{ij} to a random value in $[-\epsilon, \epsilon]$

0 Theta1 = rand(10, 11) * (2 * INIT_EPSILON) - INIT_EPSILON

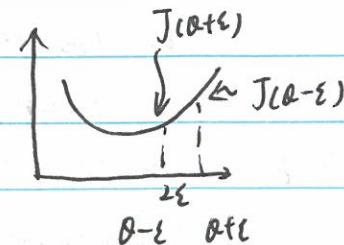
- Neural Network Learning - Putting it together

- # of input features and normalize all feature

- # of output

- # of hidden layers, have same # of hidden units in each layer
(usually the more, the better)

- Work Flow



1. Randomly initialize weight
 2. Implement forward propagation to get $h_\theta(x^i)$ for any x^i
 3. Implement to compute cost function $J(\theta)$
 4. Implement backprop to compute $\frac{\partial}{\partial \theta_j} J(\theta)$
for $j = 1: m$
- $$\Delta^L := \Delta^L + \delta^{L+1}(a^L)^T$$
5. Using gradient check to compare $\frac{\partial}{\partial \theta_j} J(\theta)$ from back propagation vs numerical estimate of gradient of $J(\theta)$
 - b. Use gradient descent or advanced optimization method with back propagation to min $J(\theta)$

Week 6

• Advice for Applying Machine Learning

- Overfitting : high variance

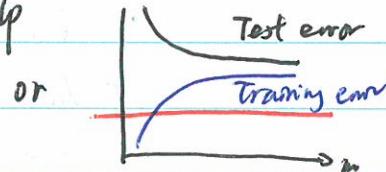
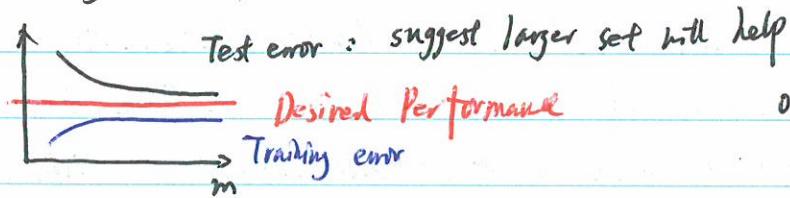
{ Too few features to classify spam (high bias)

* Diagnostic:

{ Variance : training error will be much lower than test error

{ Bias : training error will be high.

* Learning curve for high variance



* Fix high variance : more training example / smaller set of feature / $\uparrow \lambda$

{ Fix high bias : larger set of feature / email header feature / $\downarrow \lambda$

more iteration on gradient descent / Newton's method / different λ / SVM

Fix optimization algorithm :

Procedure : - * Data set is split into training set (70%) and test set (30%)

* Learning parameter θ from training set (min training error $J(\theta)$)

* Compute test set error:

$$J_{\text{test}}(\theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_i^{\text{test}}) - y_i^{\text{test}})^2$$

* Misclassification error (0/1 misclassification error)

- Model Selection and Training / Validation / Test set

* To avoid overfit problem

$$\begin{cases} \text{Split dataset: } & \begin{cases} \text{Training set (60%)}: (x_1^i, y_1^i) \dots (x_m^i, y_m^i) \\ \text{Cross validation set (20%)}: (x_{cv}^i, y_{cv}^i) \dots (x_{cv}^{m_{cv}}, y_{cv}^{m_{cv}}) \\ \text{Test set (20%)}: (x_{\text{test}}^i, y_{\text{test}}^i) \dots (x_{\text{test}}^{m_{\text{test}}}, y_{\text{test}}^{m_{\text{test}}}) \end{cases} \end{cases}$$

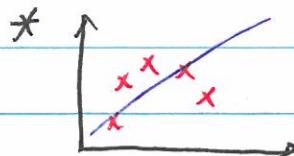
$$2. \begin{cases} \text{Training error: } J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 \\ \text{Cross Validation error: } J_{cv}(\theta) = \frac{1}{m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^i) - y_{cv}^i)^2 \end{cases}$$

$$\text{Test error: } J_{\text{test}}(\theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_i^{\text{test}}) - y_i^{\text{test}})^2$$

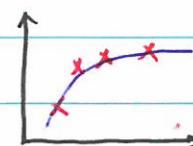
$$3. \text{pick the model: } h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

4. Report test set error: $J_{\text{test}}(\theta^*)$ and to be an optimistic estimate of generalization error.

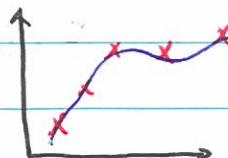
- Diagnostic bias vs Variance



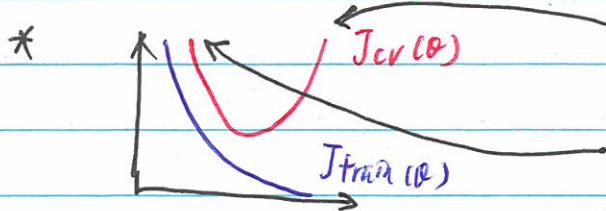
High Bias (Underfit)



Just right



High Variance (Overfit)



Variance: $J_{\text{train}}(\theta)$ is small

$J_{cv}(\theta) \gg J_{\text{train}}(\theta)$

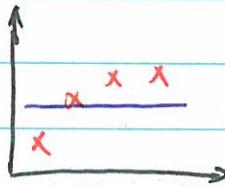
Bias: $J_{\text{train}}(\theta)$ is high

$J_{cv}(\theta) \approx J_{\text{train}}(\theta)$

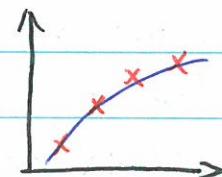
- Regularization and Bias / Variance

$$\text{Model: } h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3 + \theta_4 x_4^4 + \dots$$

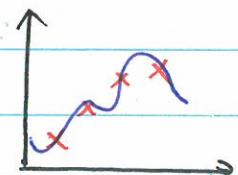
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + \lambda \sum_{j=1}^n \theta_j^2$$



large $\lambda \rightarrow$ high bias
under fit



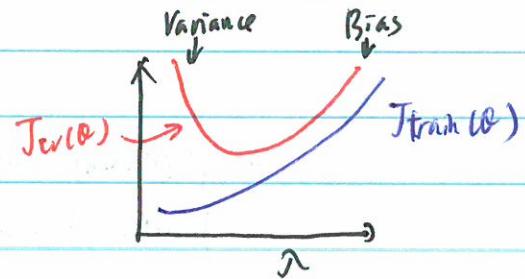
Intermediate λ
"Just right"



Small $\lambda \rightarrow$ high variance
overfit

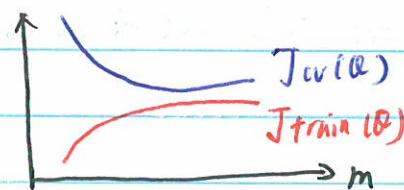
- * try different $\lambda \rightarrow \min J(\theta) \rightarrow \theta^* \Rightarrow J_{cv}(\theta^*)$ and choose min one
 \rightarrow Test error: $J_{test}(\theta^*)$

$$\begin{aligned} * J(\theta) &= \frac{1}{m} \sum_i (h_\theta(x^i) - y^i)^2 + \frac{\lambda}{m} \sum_j \theta_j^2 \\ J_{train}(\theta) &= \frac{1}{m} \sum_i (h_\theta(x^i) - y^i)^2 \\ J_{cv}(\theta) &= \frac{1}{m_{cv}} \sum_{i \in cv} (h_\theta(x_{cv}^i) - y_{cv}^i)^2 \end{aligned}$$



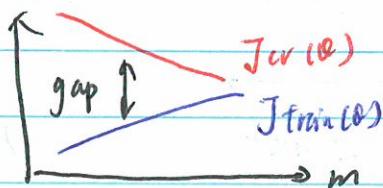
- Learning curve

- * High bias



If suffering high bias, m doesn't help

- * High variance



If suffering high variance, m is likely to help

- small neural network \rightarrow computationally cheap \rightarrow underfit
- large neural network \rightarrow computationally expensive \rightarrow overfit

Week 6

- Machine Learning System Design

- Spam classifier

- * Supervised Learning: x = features of email, y = spam (1) or not spam (0)

Features x_i : choose 100 words indicative of spam/not spam

$$x_j = \begin{cases} 1 & \text{if word } j \text{ appears} \\ 0 & \text{otherwise} \end{cases}$$

- Error Analysis

* Recommended approach:

1. Implement the algorithm and test it on cross-validation data
2. Plot learning curve to decide if more data, more features to help
3. Error analysis: manually examine the examples in cross validation set. See if you spot any systematic trend in what type of error

* Using "stemming" software: universe / university

- Error Metrics for Skewed Classes

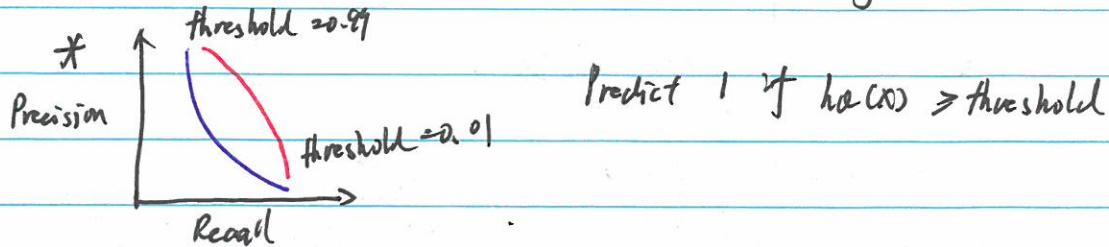
* Precision / Recall

predicted class	Actual Class		
	1	0	
1	True Positive	False Positive	
0	False Negative	True negative	

$$\text{Precision} = (\text{True Positive}) / (\text{Predicted Positive}) = (\text{True Positive}) / (\text{True Positive} + \text{False Positive})$$

$$\text{Recall} = (\text{True Positive}) / (\text{Actual Positive})$$

$$= (\text{True Positive}) / (\text{True Positive} + \text{False Negative})$$



* To predict $y=1$ only if very confident \Rightarrow higher precision, lower recall

To avoid missing too many cases of $y=1$ (avoid false negative) \Rightarrow

higher recall, lower precision

$$\text{* F}_1 \text{ Score (F score)} = 2 \cdot P \cdot R / (P + R)$$

- Large data rationale

* Using many parameters for LR or many hidden layers for NN, \Rightarrow

- \Rightarrow Lower bias and $J_{\text{train}}(\theta)$ is small
 * Using large training set \Rightarrow low variance $\Rightarrow J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$
 $J_{\text{test}}(\theta)$ is small

Week 7 Support Vector Machine

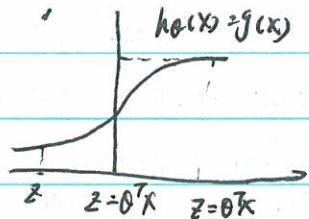
- Optimization Objective
- Alternative view of logistic regression

* $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$

If $y=1$, $h_\theta(x) \approx 1$, $\theta^T x \rightarrow \infty$

If $y=0$, $h_\theta(x) \approx 0$, $\theta^T x \rightarrow -\infty$

* $\text{Cost} = -y \log h_\theta(x) + (1-y) \log(1-h_\theta(x))$



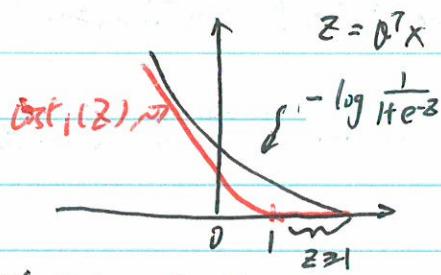
- Logistic regression:

$$\min \frac{1}{m} \left[\sum_i^m y^i (-\log h_\theta(x^i)) + (1-y^i) (-\log(1-h_\theta(x^i))) \right] + \lambda \sum_j^n \theta_j^2$$

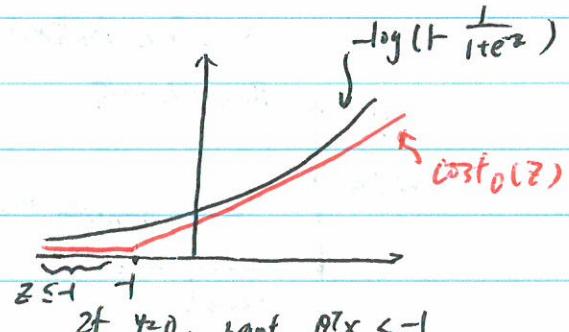
Support vector machine

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

- Large Margin Intuition



If $y=1$, want $\theta^T x \geq 1$



If $y=0$, want $\theta^T x \leq -1$

- Whenever $y^i=1$ and $\theta^T x^i \geq 1$ or $y^i=0$ and $\theta^T x^i \leq -1$

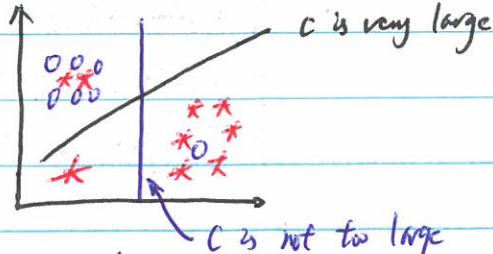
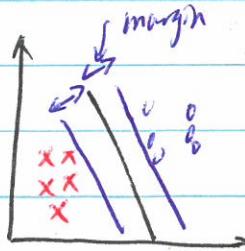
$$\text{cost} = \min \theta + \frac{1}{2} \sum_j^n \theta_j^2$$

$$\text{s.t. } \theta^T x^i \geq 1 \text{ if } y^i=1$$

$$\theta^T x^i \leq -1 \text{ if } y^i=0$$

- Decision boundary:

2f with outlier



- The mathematics behind large margin classification

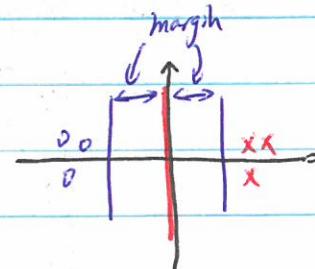
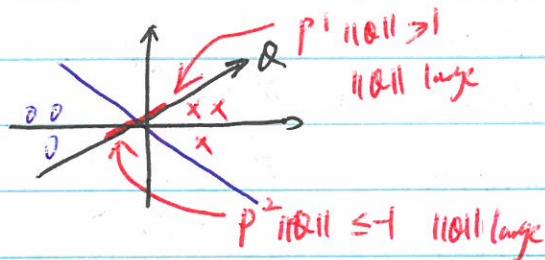
- $u = [u_1 \ u_2]$ $v = [v_1 \ v_2]$, p is the length of projection of v onto u

$$u^T v = u_1 v_1 + u_2 v_2 = p \cdot \|u\| \text{ where } \|u\| = \sqrt{u_1^2 + u_2^2}$$

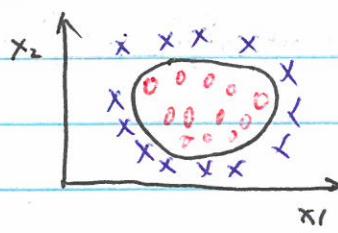
- SVM Decision boundary: $\frac{\max_j \frac{1}{2} \beta_j^2}{\frac{1}{2} \sum_j \beta_j^2} = \frac{1}{2} (\beta_1^2 + \beta_2^2) = \frac{1}{2} \|\beta\|^2$

s.t. $\beta^T x_i \geq 1$ if $y_i = 1$
 $\beta^T x_i \leq -1$ if $y_i = 0$

\Rightarrow s.t. $p_i \|\beta\| \geq 1$ if $y_i = 1$
 $p_i \|\beta\| \leq 1$ if $y_i = 0$ where p_i is the projection of x_i onto β



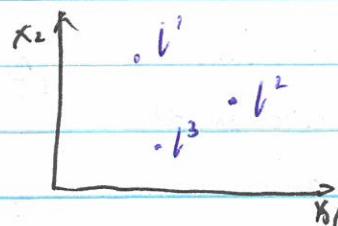
- kernels



$$h_\alpha(x) = \begin{cases} 1 & \text{if } \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 x_1^2 + \beta_5 x_2^2 \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\beta_0 + \beta_1 f_1 + \beta_2 f_2 + \beta_3 f_3 + \beta_4 f_4 + \beta_5 f_5$$

$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1 x_2, \quad f_4 = x_1^2, \quad f_5 = x_2^2$$



Given X , compute new features depending on proximity to landmark, b' , b^2 , b^3

$$\Rightarrow f_1 = \text{similarity}(X, b') = \exp\left(-\frac{\|X - b'\|^2}{2\sigma^2}\right)$$

$$f_2 = \text{similarity}(X, b^2) = \exp\left(-\frac{\|X - b^2\|^2}{2\sigma^2}\right)$$

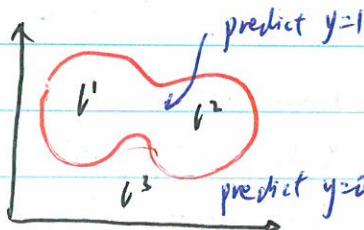
$$f_3 = \text{similarity}(X, b^3) = \exp\left(-\frac{\|X - b^3\|^2}{2\sigma^2}\right)$$

Gaussian Kernel

- $f_i = \text{similarity}(x, l^i) = \exp\left(-\frac{\|x - l^i\|^2}{2\sigma^2}\right)$

If $x \approx l^i$: $f_i \approx 1$

If x is far from l^i : $f_i \approx 0$



Predict "1" when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

- To choose l^1, l^2, l^3

* SVM with kernels

Given $(x^1, y^1) \dots (x^n, y^n)$

choose $l^1 = x^1 \dots l^m = x^m$

$$f_i = \text{similarity}(x, l^i) \Rightarrow f = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1$$

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

For training example (x^i, y^i) :

$$x^i \rightarrow f_{l^i}^i = \text{sim}(x^i, l^i)$$

$$f_{l^j}^i = \text{sim}(x^i, l^j)$$

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

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

$$x^i = 2R^{n+1}$$

$$f_i = \begin{bmatrix} f_0^i \\ f_1^i \\ \vdots \\ f_m^i \end{bmatrix} \quad 2R^n \text{ and } f_0^i = 1$$

* predict $y=1$ if $\theta^T f \geq 0$

$$\text{Training } \min_{\theta} \left(\frac{m}{2} \sum_j y_j \text{cost}(\theta^T f_j) + (1-y_j) \text{cost}_0(\theta^T f_j) + \frac{1}{2} \sum_j \theta_j^2 \right)$$

$$\text{where } \sum_j \theta_j^2 = \theta^T \theta = \|\theta\|^2$$

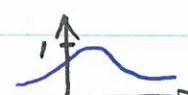
* SVM Parameters:

$C (= 1/\lambda)$: Large C : lower bias, high variance

Small C : high bias, lower variance

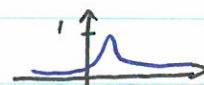
σ^2 : large σ^2 : features f_i vary more smoothly

High bias, lower variance



small σ^2 : feature f_i vary less smoothly

lower bias, higher variance



- Using SVM package (liblinear, libsvm) to solve for parameter α
- specify parameter C and choice of kernel
 - * Gaussian Kernel: $f_i = \exp\left(-\frac{\|x - l^i\|^2}{2\sigma^2}\right)$ where $l^i = x^i$
need choose σ^2
 - * Note: do perform feature scaling before Gaussian kernel
 - Other choice of kernel (need to satisfy "Mercer's Theorem")
 - * Polynomial kernel $K(x, l) = (x^T l)^d$ or $(x^T l + \text{constant})^d$
 - * String kernel, chi-square kernel, histogram intersection kernel
 - Many SVM packages have built-in multi-class classification
or use one-vs-all method
 - $n = \# \text{ of features}$, $m = \# \text{ of}$
 - * If $n \gg m$, use logistic regression or SVM w/o kernel (Predict $y=1$ if $\alpha^T x \geq 0$)
 - If n is small, m is intermediate: use SVM w Gaussian kernel
 - If $n \ll m$, add more features. use logistic regression or SVM w/o kernel
 - * Neural network work well for most, but slower to train

Week 8

Unsupervised Learning: Kmean cluster and PCA

- Unsupervised Learning = Clustering

- K-means Algorithm

* Input: $\{K \text{ (# of clusters)}$
 $\text{training set } \{x^1, x^2, \dots, x^m\}$

* Randomly initialize K clusters $\mu_1, \mu_2, \dots, \mu_K$

Repeat {

cluster assignment [for $i=1$ to m

step [$c^i = \text{index of cluster centroid closest to } x^i \left(\min_k \|x^i - \mu_k\|^2 \right)$

More centroid [for $k=1$ to K

]
 $\mu_k = \text{average (mean) of points assigned to cluster } k$

- Optimization objective

* c^i = index of cluster to which examples x^i is currently assigned

μ_k = cluster centroid k

μ_{ci} = cluster centroid of cluster to which example x^i is assigned

* Optimization objective: $x^i \rightarrow 5 \Rightarrow c^i = 5 \Rightarrow \mu_{ci} = \mu_5$

$$J(c^1 \dots c^m, \mu_1 \dots \mu_k) = \frac{1}{m} \sum_i \|x^i - \mu_{ci}\|^2$$

$$\min_{\substack{c^1 \dots c^m \\ \mu_1 \dots \mu_k}} J(c^1 \dots c^m, \mu_1 \dots \mu_k)$$

* Cluster assignment step: $\min J$ cost $c^1 \dots c^m$ (holding $\mu_1 \dots \mu_k$ fixed)

Move centroid: $\min J$ cost $\mu_1 \dots \mu_k$

- Random Initialization:

* Should have k cm

* Randomly pick k training example

Set $\mu_1 \dots \mu_k$ equal to these k examples

* For $i = 1$ to $100 \}$

Randomly initialize k-means

Run k-means, Get $c^1 \dots c^m, \mu_1 \dots \mu_k$

Compute cost function (distortion) $J(c^1 \dots c^m, \mu_1 \dots \mu_k)$

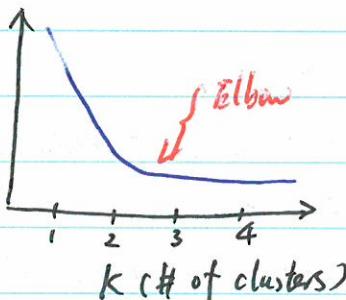
}

Pick clustering that gave lowest cost $J(c^1 \dots c^m, \mu_1 \dots \mu_k)$

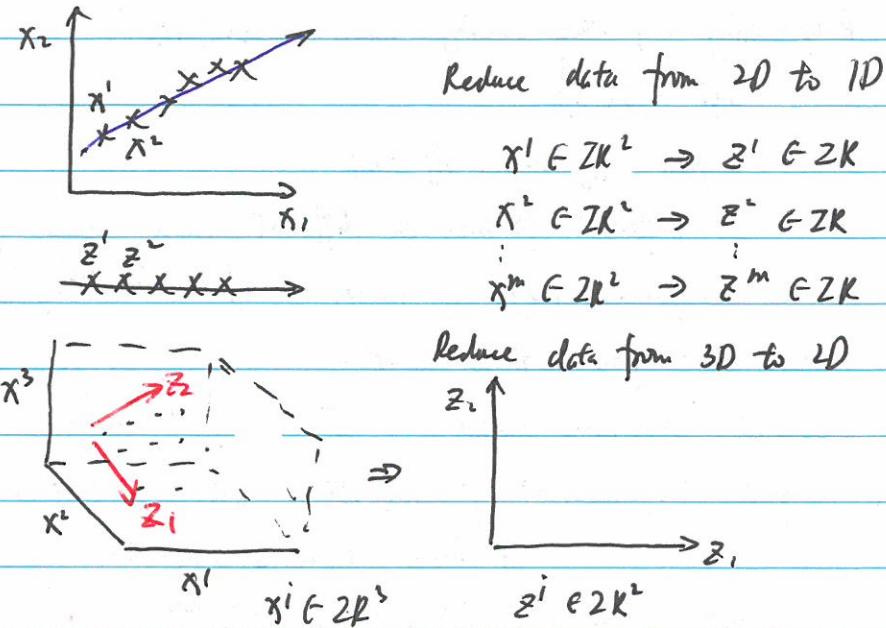
- Choosing the # of clusters

* Elbow method:

Cost Function J



- Dimensionality Reduction
- Motivation
- * Data Compression

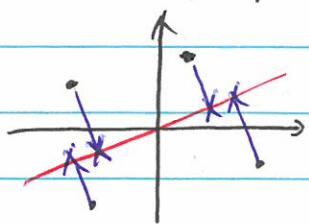


* Data Visualization

Country vs GDP, per capita GDP, Life, Poverty. Mean House Income

Reduce data from 5D to 2D (z_1, z_2), easy to plot

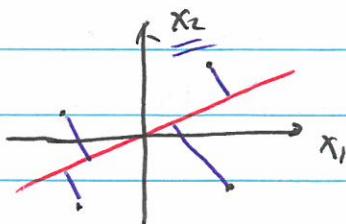
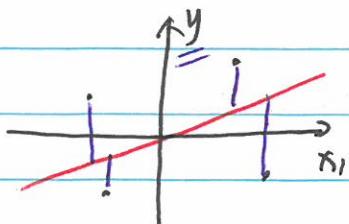
- Principal Component Analysis Problem Formulation



* Reduce from 2D to 1D: find a direction (u') onto which to project the data so as to min projection error.

* Reduce from 1D to kD: find k vector ($u^1 \dots u^k$) onto which to project the data so as to min projection error

* PCA is not linear regression



- Principal Component Analysis Algorithms

* Data Processing:

$$1. \text{ Train set} = X^1, X^2, \dots, X^m$$

2. Preprocessing (feature scaling / mean normalization)

$$\mu_j = \frac{1}{m} \sum_i x_j^i$$

$$\text{Replace } x_j^i \text{ with } x_j^i - \mu_j$$

If different features on different scales, scale features to have comparable range of values $x_j^i \in (x_j^i - \mu_j) / s_j$

* Reduce data from $n\text{-D}$ to $k\text{-D}$:

$$3. \text{ Compute "covariance matrix": } \Sigma = \frac{1}{m} \sum_i^n (x_i)(x_i)^T \Rightarrow n \times n$$

$$4. \text{ Compute "eigenvector" of matrix } \Sigma: [U, S, V] = \text{svd}(\Sigma) \quad U \in \mathbb{R}^{n \times n}$$

If singular value decomposition: $\text{eig}(\Sigma)$

$$\text{we get } U = [\underbrace{u^1, u^2, \dots, u^n}_k] \in \mathbb{R}^{n \times n} \rightarrow \text{Ureduce} = [\underbrace{u^1, u^2, \dots, u^k}_k]$$

$$5. \quad Z \in \mathbb{R}^n \Rightarrow Z \in \mathbb{R}^{n \times k}$$

$$Z^i = \underbrace{[u^1, u^2, \dots, u^k]}_{n \times k}^T x_i = \underbrace{\begin{bmatrix} -u^1 \\ -u^2 \\ \vdots \\ -u^k \end{bmatrix}}_{k \times n} x_i \quad n \times 1$$

* 1. After mean normalization (ensure $k \times 1$)

every feature has 0 mean) and optionally feature scaling

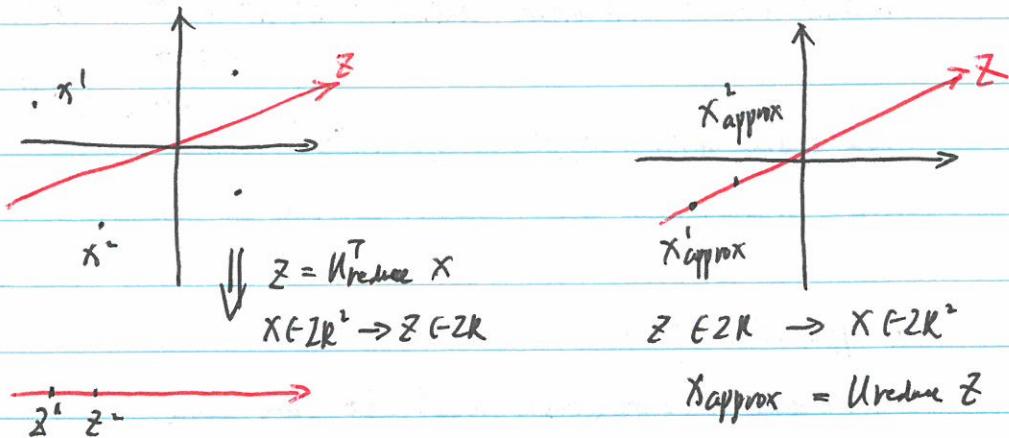
$$\Sigma (\text{sigma}) = \frac{1}{m} \sum_i^n (x_i)(x_i)^T = \frac{1}{m} * X' * X$$

$$2. [U, S, V] = \text{svd}(\text{sigma})$$

$$3. U_{\text{reduce}} = U(:, 1:k)$$

$$4. Z = U_{\text{reduce}}' * X$$

- Reconstruction from Compressed Representation



- Choosing the # of Principal Components (k)

* Average squared projection error: $\frac{1}{m} \sum_i \|x^i - x^{i\text{approx}}\|^2$

Total variation in the data: $\frac{1}{m} \sum_i \|x^i\|^2$

Choose k to be the smallest value so that

$$\frac{\frac{1}{m} \sum_i \|x^i - x^{i\text{approx}}\|^2}{\frac{1}{m} \sum_i \|x^i\|^2} \leq 0.01 \quad (1\%) \leq 99\% \text{ of variance is retained}$$

* Algorithm:

1. try PCA with $k=1$

2. Compute $Ureduce, z^1, z^2 \dots z^n, x^{i\text{approx}} \dots x^{n\text{approx}}$

3. check if $|t| \leq 0.01$

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

$$S = \begin{bmatrix} s_{11} & 0 \\ 0 & s_{22} \\ \vdots & \vdots \\ 0 & s_{nn} \end{bmatrix} \quad (\text{similar to eigenvalue})$$

For given k , $1 - \frac{\sum_i^k s_{ii}}{\sum_i^n s_{ii}} \leq 0.01$ or $\frac{\sum_i^k s_{ii}}{\sum_i^n s_{ii}} \geq 0.99$

Pick smallest value of k for $\frac{\sum_i^k s_{ii}}{\sum_i^n s_{ii}} \geq 0.99$

- Advice for Applying PCA

* Speed up Supervised Learning

Original training set: $(x^1, y^1) \dots (x^n, y^n)$

Extract input: Unbalanced dataset: $x^1, x^2 \dots x^m \in \mathbb{R}^{1000}$

↓ PCA

$$z^1, z^2 \dots z^m \in \mathbb{R}^{100}$$

New training set: $(z^1, y^1) \dots (z^m, y^m)$

Note: Mapping $x^i \rightarrow z^i$ should be only on training set. This mapping can be applied to x_{cv} and x_{test} in cross validation and test sets.

* Application of PCA

1. Reduce memory / disk need to store data

2. Speed up learning algorithm

3. Visualization

* Use z^i instead of x^i to reduce # of features to $k < n$, so to address overfitting. **Bad ideas**

use regularization instead $\min_{\theta} \frac{1}{2m} \sum_i (h_\theta(x^i) - y^i)^2 + \lambda \sum_j \theta_j^2$

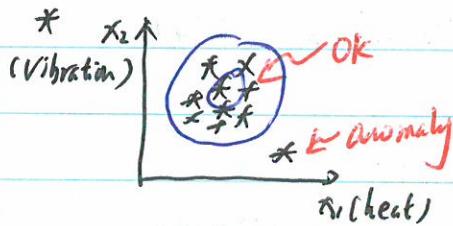
* Run Original data x^i first. if doesn't work. use PCA

Week 9

Anomaly Detection and Recommender System

- Anomaly detection

- Motivation



Aircraft engine feature: $x_1, x_2 \dots$

Dataset: $\{x^1, x^2 \dots x^m\}$

New engine: x_{test}

Model: $P(x)$

$P(x_{\text{test}}) < \epsilon \Rightarrow \text{Anomaly}; P(x_{\text{test}}) \geq \epsilon \Rightarrow \text{OK}$

* Application: fraud detection / Manufacture / Monitoring computer in data center

- Gaussian distribution:

* $x \sim N(\mu, \sigma^2)$

$$P(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad \left\{ \begin{array}{l} \mu = \frac{1}{m} \sum_i^m x^i \\ \sigma^2 = \frac{1}{m} \sum_i^m (x^i - \mu)^2 \end{array} \right.$$



$G \downarrow \Rightarrow \text{sharp.}$

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

- Algorithm

* Density estimation

Training set: $\{x^1, \dots, x^m\}$

$$x_i \sim N(\mu_i, \sigma_i^2) \quad x_1 \sim N(\mu_1, \sigma_1^2) \quad x_2 \sim N(\mu_2, \sigma_2^2) \quad x_3 \sim N(\mu_3, \sigma_3^2)$$

$$\begin{aligned} P(x) &= P(x_1, \mu_1, \sigma_1^2) P(x_2, \mu_2, \sigma_2^2) \cdots P(x_n, \mu_n, \sigma_n^2) \\ &= \prod_{j=1}^n P(x_j, \mu_j, \sigma_j^2) \end{aligned}$$

* Algorithm

1. Choose features x_i that might be indicative of anomalous examples

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

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_{ij}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_{ij} - \mu_j)^2$$

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

3. Given new example x , compute $P(x)$:

$$P(x) = \prod_{j=1}^n P(x_j, \mu_j, \sigma_j^2)$$

Anomaly if $P(x) < \epsilon$

- Anomaly Detection System

* 1. Choosing feature x_i

2. Label anomalies vs non-anomalies (1 vs 0)

3. Training set: $x^1 \dots x^m$ 60% good, 0% bad

Cross validation set: $(x_{cv}^1, y_{cv}^1) \dots (x_{cv}^{new}, y_{cv}^{new})$ 20% good, 50% bad

Test set: $(x_{test}^1, y_{test}^1) \dots (x_{test}^{new}, y_{test}^{new})$ 20% good, 50% bad

4. Fit model $P(x)$ on training set $x^1 \dots x^m$ y

5. On cross validation/test, predict

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

Possible metrics: { True positive, false positive, false negative, true negative
Precision/Recall
F1-score

Note: can use cross-validation set to choose ϵ

- Anomaly detection vs supervised Learning

- * small # of positive ($y=1$)
- * large # of negative ($y=0$)
- * Many different type of anomalies
- * future anomalies may look nothing like any previous ones
- * Fraud detection / Manufacture / Monitor machine

Large # of positive and negative

Algorithm can get a sense of what positive examples are like, future is similar to old ones

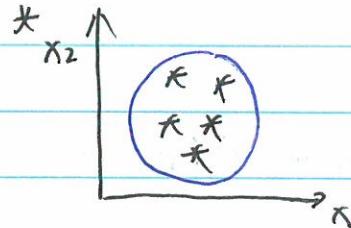
Email spam / Cancer / Weather prediction

- Choosing What Features to Use

- * Non-gaussian features: hist / log
- * Err analysis:

$p(x)$ larger for normal, $p(x)$ small for anomalies

- Multivariate Gaussian Distribution



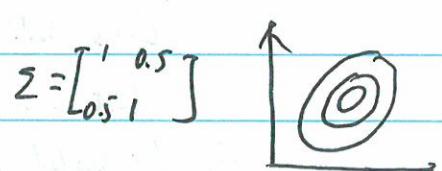
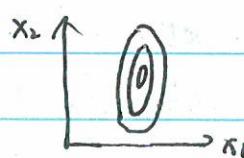
Don't model $p(x_1), p(x_2)$ separately. Model $p(x)$ all in one

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

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

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

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



1. Parameter: μ, Σ

$$p(x, \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)) \quad (\text{#})$$

2. Parameter fitting:

$$\mu = \frac{1}{m} \sum_i^m x^i \quad \Sigma = \frac{1}{m} \sum_i^m (x^i - \mu)(x^i - \mu)^T$$

3 given a x_{test} to (#), flag anomaly if $p(x) < \epsilon$

* For $\Sigma = \begin{bmatrix} \sigma_i^2 & 0 \\ 0 & \sigma_n^2 \end{bmatrix}$

* Original Model vs Multivariate Gaussian
 $\prod P(x_j | \mu_j, \sigma_j^2)$ $P(x, \mu, \Sigma)$

Independent feature

Computationally cheap

O(k) even if m small

Captures correlation between features

Computationally expensive

Must have $m > n$ or else Σ is non-invertible ($m \geq n$)

- Recommender Systems

- Problem formulation

* Predict movie ratings from 0 to 5 stars

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 (only if $r(i, j) = 1$)

- Content-based Recommendation

* For each user j , learn a parameter $\theta^j \in \mathbb{R}^{n+1}$

Predict user j as rating movie i with $(\theta^j)^T x^i$ stars

* Problem formulation

$r(i, j) = 1$ if user j has rated movie i

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

$\theta^{(i)}$ = parameter vector for user j ($\theta^j \in \mathbb{R}^{n+1}$)

x^i = feature vector for movie i

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

m_j = # of movies rated by user j

Optimization objective : to learn θ^j

$$\min_{\theta^j} \frac{1}{2} \sum [(\theta^j)^T x^i - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{k=1}^n (\theta_k^j)^2$$

to learn $\alpha^1, \alpha^2, \dots, \alpha^{n_u}$:

$$\min_{\alpha^1, \dots, \alpha^{n_u}} \frac{1}{2} \sum_{j=1}^{n_u} \sum [(\alpha^j)^T x^j - y^{i,j}]^2 + \lambda \sum_{k=1}^{n_u} \sum (\alpha_k^j)^2$$

Gradient descent update: $J(\alpha^1, \dots, \alpha^{n_u})$

$$\alpha_k^j = \alpha_k^j - \alpha \sum [(\alpha^j)^T x^j - y^{i,j}] x_k^j \quad (\text{for } k \geq 0)$$

$$\alpha_k^j = \alpha_k^j - \alpha \left\{ \sum [(\alpha^j)^T x^j - y^{i,j}] + \lambda \alpha_k^j \right\} \quad (\text{for } k \neq 0)$$

- Collaborative Filtering

* Given $\alpha^1 \dots \alpha^{n_u}$, to learn x^i :

$$\min_x \frac{1}{2} \sum [(\alpha^j)^T x^i - y^{i,j}]^2 + \lambda/2 \sum (x_k^i)^2$$

Given $\alpha^1 \dots \alpha^{n_u}$ to learn $x^1 \dots x^{n_m}$

$$\min_x \frac{1}{2} \sum \sum [(\alpha^j)^T x^i - y^{i,j}]^2 + \lambda/2 \sum \sum (x_k^i)^2 \quad (1)$$

* Given $x^1 \dots x^{n_m}$ can estimate $\alpha^1 \dots \alpha^{n_u}$

$$\min_{\alpha} \frac{1}{2} \sum \sum [(\alpha^j)^T x^i - y^{i,j}]^2 + \lambda/2 \sum \sum (\alpha_k^j)^2$$

Given $\alpha^1 \dots \alpha^{n_u}$ can estimate $x^1 \dots x^{n_m}$ (2)

$$\min_x \frac{1}{2} \sum \sum [(\alpha^j)^T x^i - y^{i,j}]^2 + \lambda/2 \sum \sum (x_k^i)^2$$

Min $x^1 \dots x^{n_m}$ and $\alpha^1 \dots \alpha^{n_u}$ simultaneously:

$$J = (1) + (2)$$

* Algorithm:

1. Initialize $x^1 \dots x^{n_m}$, $\alpha^1 \dots \alpha^{n_u}$ to small random values

2. Min J using gradient descent.

$$x_k^i = x_k^i - \alpha \left(\sum [(\alpha^j)^T x^i - y^{i,j}] \alpha_k^j + \lambda x_k^i \right)$$

$$\alpha_k^j = \alpha_k^j - \alpha \left(\sum [(\alpha^j)^T x^i - y^{i,j}] x_k^i + \lambda \alpha_k^j \right)$$

3. For a user with α and a movie with x , predict star rating of $\alpha^T x$

- Vectorization: Low Rank Matrix Factorization

$$x = \begin{bmatrix} -(x^1)^T \\ \vdots \\ -(x^{n_m})^T \end{bmatrix} \quad \alpha = \begin{bmatrix} -(\alpha^1)^T \\ \vdots \\ -(\alpha^{n_u})^T \end{bmatrix}$$

$$Y = (\theta^j)^T \cdot X^i = \begin{bmatrix} (\theta^j)^T x^i & \dots & (\theta^{(m)})^T x^i \\ (\theta^j)^T x^1 & \dots & (\theta^{(m)})^T x^1 \\ \vdots & \ddots & \vdots \\ (\theta^j)^T x^{(m)} & \dots & (\theta^{(m)})^T x^{(m)} \end{bmatrix}$$

* How to find movie j related to movie i

smallest $\|x^i - x^j\| \Rightarrow$ movie i and j are similar

- Implementational Detail: Mean Normalization

* For user j , on movie i predict: $(\theta^j)^T x^i + b_i$

Week 10

Large-scale Machine Learning

- Learning with large datasets

$$\theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i) x_j^i \quad \text{where } m = 1000000$$

- Stochastic Gradient Descent

- Linear regression with gradient descent

$$h_\theta(x) = \sum_{j=0}^n \theta_j x_j \quad J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2$$

Repeat {

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

Batch Gradient Descent

all m examples

- Stochastic Gradient Descent

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

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

* 1. Randomly shuffle dataset

2. Repeat { for $i=1 \dots m$

$$\theta_j = \theta_j - \alpha [h_\theta(x^i) - y^i] \cdot x_j^i \quad \text{use 1 example}$$

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

- Mini-batch Gradient Descent

- Batch gradient descent: use all m examples in each iteration

- Stochastic gradient descent: use 1 examples in each iteration

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

* Get $b=10$ examples $(x^i, y^i) \sim (x^{i+9}, y^{i+9})$

$$\theta_j = \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (\hat{h}_\theta(x^k) - y^k) x_j^k$$

or $i = 1, 11, 21, 31, \dots, 91$

- Stochastic gradient descent convergence

- Batch gradient descent : plot $J(\theta)$ vs # of iteration

- Stochastic gradient descent :

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

- * Every 100 iteration, plot $\text{cost}(\theta, (x^i, y^i))$ averaged over last 100 example

- * Learning rate α is held constant. Can slowly decrease α over time if we want θ to converge ($\alpha = \text{const} / (\text{iteration number} + \text{const})$)

Note: as iteration number \uparrow , $\alpha \rightarrow 0$

- Online Learning

- Shipping service : feature x , use service ($y=1$), not ($y=0$)

Learn $P(y=1|x; \theta)$ to optimize price

- Product search : feature x , click link ($y=1$) not ($y=0$)

Learn $P(y=1|x; \theta)$

- Map reduce and Data Parallelism

Batch gradient descent : $\theta_j = \theta_j - \alpha \frac{1}{400} \sum_{i=1}^{400} (\hat{h}_\theta(x^i) - y^i) x_j^i$

Machine 1: $\text{temp}_j^1 = \sum_{i=1}^{100} (\hat{h}_\theta(x^i) - y^i) x_j^i$

Machine 2: $\text{temp}_j^2 = \sum_{i=101}^{200} (\hat{h}_\theta(x^i) - y^i) x_j^i$

Machine 3: $\text{temp}_j^3 = \sum_{i=201}^{300} (\hat{h}_\theta(x^i) - y^i) x_j^i$

Machine 4: $\text{temp}_j^4 = \sum_{i=301}^{400} (\hat{h}_\theta(x^i) - y^i) x_j^i$

$$\theta_j = \theta_j - \alpha \frac{1}{400} \sum_{i=1}^{400} \text{temp}_j^i$$

Multi-core machine : Core 1, Core 2, Core 3, Core 4