



Machine Learning : Andrew NG

Machine learning

given data,

- Arthur Samuel (1958)

- Tom Mitchell: E/T/P

Supervised

- "right answers" given
- types of supervised learning problems

Regression
(continuous-valued
attribute)

Classification
(discrete)

Unsupervised

- little/no info abt
"right answers"
- derive structure from data, where we
don't know effect of variables
- no feedback based on prediction results

Clustering

Non-clustering

Cocktail party problem algorithm.
Structure in chaotic env.

Supervised learning problem: given a training set, to learn a function $h: x \rightarrow y$, so that $h(x)$ is a "good" predictor for corresponding value of y

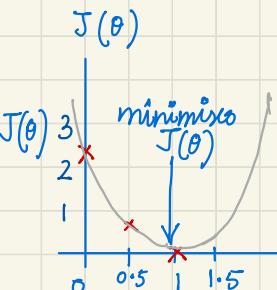
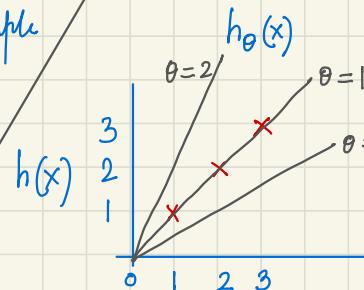
Training set: m, x, y

(x, y) : one training example

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

Learning algorithm

$x \rightarrow h \rightarrow$ predicted y
hypothesis



Problem: Univariate Linear Regression

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

Parameters: θ_0, θ_1

Cost function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_\theta(x^i) - y^i]^2$

→ for fixed θ , f^n of x → funcⁿ of θ

Goal: minimize $J(\theta_0, \theta_1)$

Assignment vs Truth assertion	
$a := b$	$a = b$
$a := a + 1$	$a = a + 1 \times$

→ Repeat until convergence?

$$\left. \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \\ \quad \quad \quad \text{for } j = 0 \text{ and } j = 1 \end{array} \right\}$$

learning rate
simultaneous update
adjust: failure / long time to converge
derivatives

Linear Algebra Review:

→ Matrices and Vectors

$$\begin{matrix} \times \\ m \text{ rows} \\ n \text{ columns} \\ \vdots \\ A_{ij} \\ \vdots \\ n = j \end{matrix}$$

→ **Matrix**: 2-dimensional array of numbers, where elements $A_{ij} \in \mathbb{R}^{m \times n}$

→ **Vector**: $n \times 1$ matrix : $\begin{bmatrix} \vdots \\ y_i \\ \vdots \\ y_j \end{bmatrix} \quad y_i \in \mathbb{R}^n$
 $\equiv n\text{-dim vector} \quad n=1$

→ **Scalar**: object is a single value

→ Matrix addition and subtraction

↳ same dimension / element-wise / BODMAS

→ Matrix Multiplication

→ Matrix - vector multiplication

→ matrix - matrix multiplication

→ Properties:

• $AB \neq BA$ (non-commutative)

• $(AB)C \neq A(BC)$ (non-associative)

• $A \cdot I = I \cdot A = A$ (Identity matrix)

$m \times n \quad n \times m \quad m \times n$ dimension implicit

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

→ Inverse and Transpose

• **Matrix Inverse**: If $A = m \times m$ and has inverse

$$AA^{-1} = A^{-1}A = I$$

• **Transpose**: $A \rightarrow A^T$

$$A_{ij} = A_{ji}^T$$

Singular matrix



Multi-variate linear regression

Multiple features: n , $x_j^{(i)}$ (i)th training example
 $\#$ features j th feature
 m : # training examples

$$\text{Hypothesis: } h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 \\ = \sum_{i=1}^n \theta_i x_i \quad \text{Define } x_0 = 1 = \theta^T x$$

One Hypothesis:

House Sizes

$$\rightarrow 2104$$

$$\rightarrow 1416$$

$$\rightarrow 1534$$

$$\rightarrow 852$$

Hypothesis

$$h_\theta(x) = -40 + 0.25x$$

Matrix

$$\begin{bmatrix} 2104 \\ 1416 \\ 1534 \\ 852 \end{bmatrix}$$

Vector

$$\begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$$

$$4 \times 2 \times 1 = 4 \times 1$$

$$\text{Data matrix} \times \text{parameters} = \text{Prediction}$$

Computing hypothesis:

House sizes:

$$\begin{cases} \dots \\ \dots \\ \dots \\ \dots \end{cases}$$

Hypothesis:

$$1. h_\theta(x) = \dots$$

$$2. h_\theta(x) = \dots$$

$$3. h_\theta(x) = \dots$$

Matrix

$$\begin{bmatrix} \dots \\ \dots \\ \dots \\ \dots \end{bmatrix}$$

Matrix

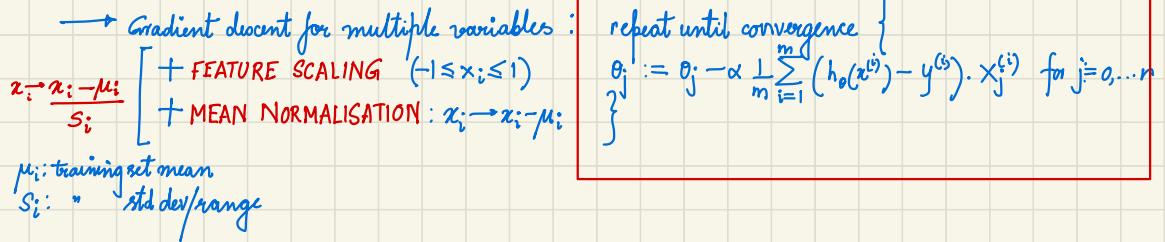
$$\begin{bmatrix} \dots \\ \dots \\ \dots \\ \dots \end{bmatrix} = \text{Prediction}$$

Univariate \rightarrow multivariate

x : scalar

x : vector

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$



- + DEBUGGING : $J(\theta)$ vs #iterations → slow convergence, α too small
- + Automatic convergence test → $J(\theta) \neq$, non-convergence, α too large
- + LEARNING RATE Try $0.001, 0.003, 0.01, 0.03, 0.1, \dots$

→ Linear Regression vs Polynomial Regression

Choice depends on specific training set example

converting polynomial → linear problem

$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

Making new features by combining

$$\begin{cases} x_1 = x & \text{Ranges change} \\ x_2 = x^2 \\ x_3 = x^3 \end{cases}$$

Minimisation algorithms : for $J(\theta)$

→ Gradient Descent

→ Normal Equation :

$$\theta \in \mathbb{R}^{n+1}, J(\theta_0, \theta_1, \dots, \theta_n)$$

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

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

$$\frac{\partial J(\theta)}{\partial \theta_j} = \dots \stackrel{\text{set}}{=} 0 \quad \forall j$$

Solve for $\theta_0, \theta_1, \dots, \theta_n$



$$\theta = (X^T X)^{-1} X^T y$$

per training eq

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

X : Design matrix = $\begin{bmatrix} \dots (x^{(1)})^T \dots \\ \dots (x^{(2)})^T \dots \\ \dots (x^{(m)})^T \dots \\ m \times (n+1) \end{bmatrix}$

Gradient Descent

→ iterative minimisation of $J(\theta)$

$$\frac{\partial J(\theta)}{\partial \theta_j}$$

Normal Eqⁿ

→ explicitly & analytically compute $\frac{\partial J(\theta)}{\partial \theta_j} = 0$ to

solve for θ_j

Use Normal Eqⁿ

$$\theta = (X^T X)^{-1} X^T y$$

→ choose α

→ needs many iters

→ $O(kn^2)$

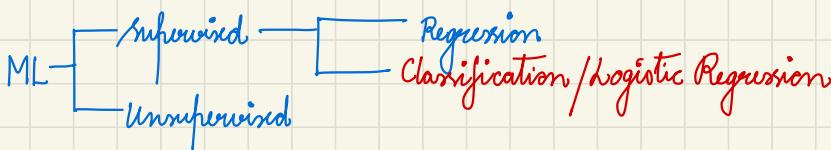
→ works well for large N

→ no need for α

→ no need for iters

→ $O(n^3)$, need to compute $(X^T X)^{-1}$

→ computationally limited



CLASSIFICATION / Logistic Regression

- Binary classification: $y \in \{0, 1\} \equiv \{-, +\}$ \rightarrow Want: $0 \leq h_\theta(x) \leq 1$

Hypothesis Representation
training example label

Model: $h_\theta(x) = g(\theta^T x)$

$$g(z) = \frac{1}{1+e^{-z}}$$

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

Interpretation: estimated probability that $y=1$ on input x , parameterized by θ

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



$$\begin{aligned} h_\theta(x) \geq 0.5 &\rightarrow y=1 \leftarrow \theta^T x \geq 0 \\ h_\theta(x) < 0.5 &\rightarrow y=0 \leftarrow \theta^T x < 0 \end{aligned}$$

$$\begin{aligned} P(y=0|x; \theta) + P(y=1|x; \theta) &= 1 \\ P(y=0|x; \theta) &= 1 - P(y=1|x; \theta) \end{aligned}$$

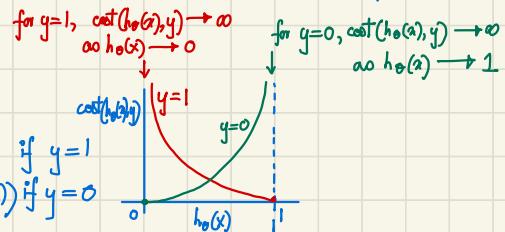
Decision boundary: $h_\theta(x) = 0.5 \Leftrightarrow \theta^T x = 0$

↳ property of hypothesis function
↳ linear/polynomial

penalty when prediction with abs. certainty turns out to be wrong!

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

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



$$\begin{aligned} J(\theta) &= \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right] \\ &= \frac{1}{m} (-y^T \log(h) - (1-y)^T \log(1-h)) \end{aligned}$$

→ Minimize $J(\theta)$ using GRADIENT DESCENT

Repeat {

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

(simultaneously update all θ_j)

looks identical to linear regression but

$$h_\theta(x) = \theta^T x$$

vs

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

$$\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \vec{y}) \quad [\text{vectorized implementation}]$$

{ [Can also use ready-to-use advanced algorithms]

- Multi-class classifications: One vs All (one-vs-rest)

$$y \in \{0, 1, \dots, n\}$$

→ Train logistic regression classifier/hypothesis

$$h_{\theta}^{(i)}(x) = p(y=i|x, \theta) \text{ for each class } i$$

→ For new input x , to make prediction, pick class i , which maximises: $\max_i h_{\theta}^{(i)}(x)$

Underfitting/ high bias vs Overfitting/ high variance

- hypothesis doesn't map well to data trend
- too simple
- too few features
- fits data but doesn't generalise to predict new data
- too complicated
- too many

Solutions:

1. ↓ # features
 - manual
 - model selection

2. Regularisation

- smaller values for params: $\theta_0, \dots, \theta_n$
- "Simpler" hypo
- Less prone to overfitting

→ Takes care of non-invertibility

Non-invertibility:

if $m < n$
 $X^T X$ non-inv

but $X^T X + \lambda I$ is inv.

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

Regularisation parameter

- Algorithms

→ Gradient Descent:

Repeat {

$$\theta_0 := \theta_0 - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \left(1 - \frac{\alpha}{m}\lambda\right)\theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \forall j \neq 0$$

→ Normal Equation

$$\theta = (X^T X + \lambda I)^{-1} X^T y$$

where $I = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}_{(n+1) \times (n+1)}$

LOGISTIC REGRESSION

- Regularised Cost function

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

- Algorithms:

→ Gradient descent: same as linear regression, but w/ different hypothesis

→ advanced optimisation algorithms

NEURAL NETWORKS

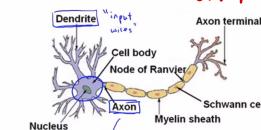
Learning algorithms:

- Linear Regression [Regression, e.g., housing prices]
- Logistic Regression [Classification, e.g., tumor prediction]
- Neural Networks [Non-linear hypothesis]

} not scalable to non-linear hypothesis with many features

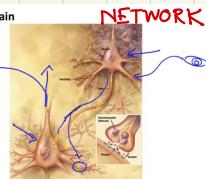
Approximation of brain's learning algorithm

Neuron in the brain



UNIT

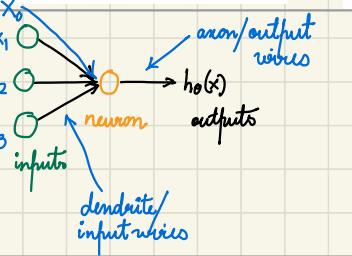
Neurons in the brain



"bias unit" = 1

Single Neuron

Neural Model: Logistic Model

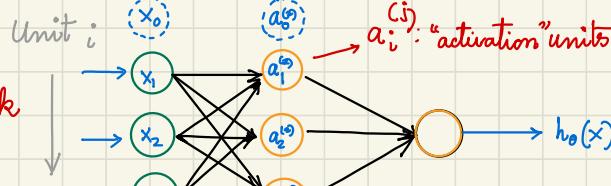


$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}, \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

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

Sigmoid (logistic) activation function

"weights" / "params"



Simplistic Representation

$$[x_1 \ x_2 \ \dots \ x_n] \rightarrow [\quad] \rightarrow h_\theta(x)$$

$$a^{(0)} = x \quad (\text{add } a_0^{(0)})$$

$$z^{(0)} = \Theta^{(0)} a^{(0)} \quad a^{(0)} = g(z^{(0)})$$

(add $a_0^{(0)}$)

$$z^{(1)} = \Theta^{(1)} a^{(0)} \quad a^{(1)} = g(z^{(1)}) = h_\theta(x)$$

Layer j → Layer 1 / Input layer

Layer 3 / Hidden layers

Layer 3 / Output layer

Forward propagation:
Vector implementation

$\Theta^{(j)}$: matrix of weights controlling function mapping from layer j → j+1

Network: layer j → s_j units

$$\Theta^{(j)} \in \mathbb{R}^{s_{j+1} \times (s_j + 1)}$$

Generalizing:

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad z^{(j)} = \begin{bmatrix} z_0^{(j)} \\ z_1^{(j)} \\ \vdots \\ z_n^{(j)} \end{bmatrix}$$

Setting $x = a^{(j)}$

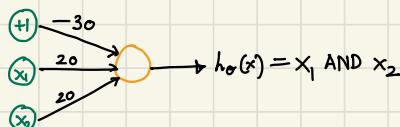
$$z^{(j)} = \oplus^{(j-1)} a^{(j-1)}$$

$s_j \times 1 \quad s_j \times (n+1) \quad (n+1) \times 1$

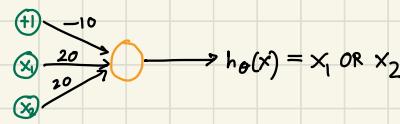
Examples and Intuitions:

→ Construct fundamental operations
in computers (AND/OR/XOR/XNOR)
using neural networks

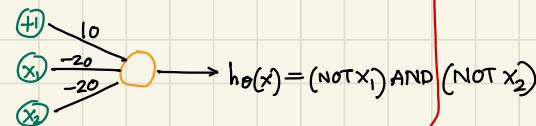
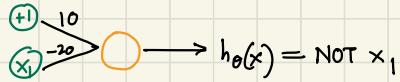
AND



OR



NEGATION



$$\text{Layer } j \Rightarrow a^{(j)} = g(z^{(j)}) \text{ [element-wise]}$$

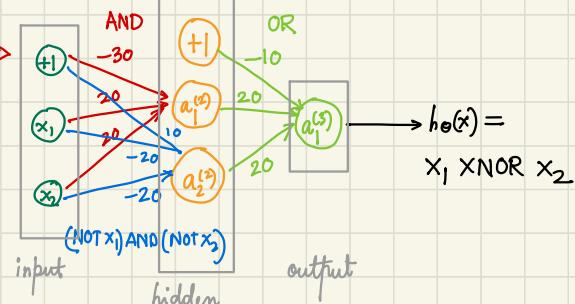
Add bias unit to layer j : $a_0^{(j)}$
 $a^{(j)} \in \mathbb{R}^{s_j} \rightarrow \mathbb{R}^{s_{j+1}}$

$$\text{Layer } j+1 \Rightarrow z^{(j+1)} = \oplus^{(j+1)} a^{(j)}$$

$1 \times 1 \quad 1 \times (s_j + 1) \quad (s_j + 1) \times 1$

$$\text{Final step: } h_\theta(x) = a^{(j+1)} = g(z^{(j)})$$

Same as Logistic Regression



Multi-class Classification

$$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ \vdots \\ a_m^{(2)} \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(3)} \\ a_1^{(3)} \\ \vdots \\ a_m^{(3)} \end{bmatrix} \rightarrow \dots \rightarrow \begin{bmatrix} h_\theta(x)_1 \\ h_\theta(x)_2 \\ \vdots \\ h_\theta(x)_m \end{bmatrix}$$



COST FUNCTION [Neural Networks]

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log h_\theta(x_k^{(i)}) + (1-y_k^{(i)}) \log (1-h_\theta(x_k^{(i)})) \right] + \frac{\lambda}{m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\Theta_{j,i}^{(l)})^2$$

L: total # layers in the network

S_l : # units (not counting bias unit) in layer l

K: # output units/classes

$h_\theta(x)_k$: hypothesis resulting in kth output

(i,j,k,l) : unit in lth layer, unit in (l+1)th layer, output, layer.
 $\Theta_{j,i}^{(l)}$ $\in \mathbb{R}$

Gradient Computation: Backpropagation

= Neural network terminology for minimising cost function

BP in NNs = Gradient descent in linear/logistic regression

Goal: $\min_{\Theta} J(\Theta) \rightarrow$ compute $\frac{\partial}{\partial \Theta_{j,i}^{(l)}} J(\Theta) \leftarrow$

Algorithms:

- gradient descent
 - ↓ advanced (fminc, fminng)
 - ↓ BP
- normal eq's

→ Backward propagation algorithm

→ Given training set $\{(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})\}$

- Set $\Delta_{ij}^{(l)} := 0 \quad \forall i,j,l$ (initializing)

$$\frac{\partial}{\partial \Theta_{j,i}^{(l)}} J(\Theta) = D_{j,i}^{(l)}$$

→ For training example t = 1 → m,

1. Set $a^{(0)}_j := x^{(t)}$

2. Perform FP to compute $a^{(l)}$ for $l = 2, 3, \dots, L$

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

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

3. Using $y^{(t)}$, compute $s^{(l)} = a^{(l)} - y^{(t)}$ (error values)

4. Perform BP to compute $s^{(L-1)}, s^{(L-2)}, \dots, s^{(2)}$, where:

$$s^{(l)} = ((\Theta^{(l)})^T s^{(l+1)}) * \underbrace{a^{(l)}}_{|||} * \underbrace{(1 - a^{(l)})}_{g'(z^{(l)})}$$

5. Update Δ matrix:

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_i^{(l)} s_i^{(l+1)}$$

$$\text{or, } \Delta^{(l)} := \Delta^{(l)} + s^{(l+1)} (a^{(l)})^T$$

$$\begin{aligned} \bullet D_{ij}^{(l)} &:= \frac{1}{m} (\Delta_{ij}^{(l)} + \lambda \Theta_{j,i}^{(l)}) \text{ if } j \neq 0 \\ \bullet D_{0j}^{(l)} &:= \frac{1}{m} \Delta_{0j}^{(l)} \quad j = 0 \end{aligned}$$

Learning Algorithm

Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.

→ Unroll to get initialmeta to pass to

→ fmminunc (costFunction, initialmeta, options)

function [jval, gradientVec] = costFunction(thetaVec)

From thetaVec, get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.

Use forward prop/back prop to compute $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\Theta)$.
 Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradientVec.

One effective strategy for choosing λ is to base it on the number of units in the network. A good choice of λ is $\lambda = \sqrt{\frac{\lambda}{n_{\text{units}}}}$, where $n_{\text{units}} = 8$ and $L_{\text{out}} = 8$ are the number of units in the layers adjacent to $\Theta^{(l)}$.

MACHINE-LEARNING DIAGNOSTICS

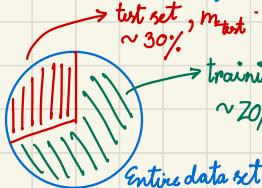
Basic trouble-shooting techniques:

- more training examples → fixes high variance
- buying smaller set of features (counter overfitting) → fixes high variance
- buying additional features (counter underfitting) → fixes high bias
- buying polynomial features (include nonlinearities) → fixes high bias
- ↑ or ↓ λ →

→ Evaluating Hypothesis :

→ learn Θ and $\min J_{\text{train}}(\Theta)$ using training set

→ compute test set error $J_{\text{test}}(\Theta)$



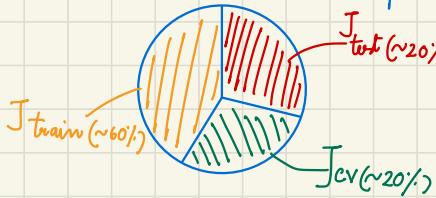
$$\text{err}(h_{\Theta}(x), y) = \begin{cases} 1 & \text{if } h_{\Theta}(x) \geq 0.5 \text{ and } y = 0 \\ 0 & \text{if } h_{\Theta}(x) < 0.5 \text{ and } y = 1 \end{cases}$$

[linear regression]

[classification]

proportion of cases misclassified

→ Model Selection : different polynomial degrees (d)



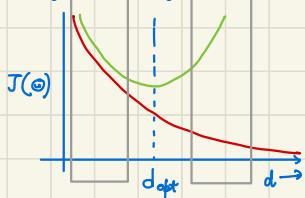
- optimise Θ for different d for training set
- use CV to find d with least J_{cv}
- estimate generalisation error using $J_{\text{test}}(\Theta^{(d)})$

↓
polynomial is not trained using test set

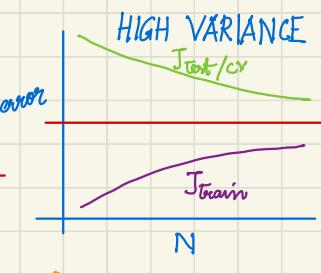
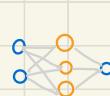
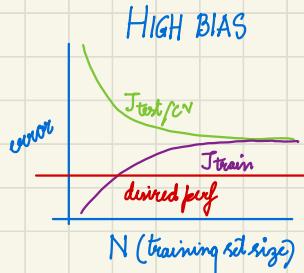
→ Diagnosing Bias vs Variance

underfitting / high λ overfitting / low λ

① Polynomial degree, d



LEARNING CURVES



→ ② Regularisation

- create λ list, $\lambda \in 0, 0.1, 0.2, 0.4, \dots, 10^{-24}$
- list of models with different λ
- Iterate through λ , and $\forall j$, through J , $J_{\text{train}}(\theta)$
- Compute $J_{\text{cv}}(\theta_{\text{train}})$ with $\lambda=0$
- Select (λ, j) combo for least J_{cv}
 use λ
- generalise to J_{test}

Computationally cheap but prone to underfitting

expensive but prone to overfitting

Use regularisation to counter

BUILDING A SPAM CLASSIFIER

→ Prioritising what to work on

→ Recommended approach

↳ quick and dirty implementation on CV set + learning curves
(pre-optimisation)

↳ Error analysis (manually examine)

↳ common types of misclassified examples (helps what to focus on first), e.g. pharma (12), replica (4), steal password (53)

↳ cues what might help prevent these misclassifications, e.g., misspellings (5), unusual punctuation (32)

↳ Numerical evaluation, e.g., "Porter stemmer" / stemming

SKEWED DATA

$y=1$: rare case

		Actual
Predicted	1	True Pos
	0	False Pos False Neg

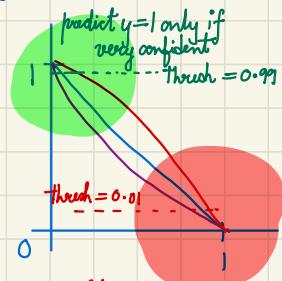
$$\rightarrow \text{Precision} = \frac{\text{True Pos}}{\#\text{Predicted pos}} = \frac{\text{True Pos}}{\text{True Pos} + \text{False Pos}}$$

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

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

$$F_1 \text{ Score} : \frac{2 \text{PR}}{\text{P} + \text{R}}$$

$$F\text{-Score} = \begin{cases} 0 & \text{if } \text{P}=0 \text{ or } \text{R}=0 \\ 1 & \text{if } \text{P}=1 \text{ & } \text{R}=1 \end{cases}$$



avoid too many false negatives

LARGE DATA SETS

→ Rationale: for low bias algorithm, larger dataset = less chances of overfitting

↳ "enough" info: tip: given input x , will human expert be able to confidently predict y ?

Optimization Objective

The **Support Vector Machine (SVM)** is yet another type of *supervised* machine learning algorithm. It is sometimes cleaner and more powerful.

Recall that in logistic regression, we use the following rules:

if $y=1$, then $h_\theta(x) \approx 1$ and $\Theta^T x \gg 0$

if $y=0$, then $h_\theta(x) \approx 0$ and $\Theta^T x \ll 0$

Recall the cost function for (unregularized) logistic regression:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m -y^{(i)} \log(h_\theta(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \\ = \frac{1}{m} \sum_{i=1}^m -y^{(i)} \log\left(\frac{1}{1 + e^{-\theta^T x^{(i)}}}\right) - (1 - y^{(i)}) \log\left(1 - \frac{1}{1 + e^{-\theta^T x^{(i)}}}\right)$$

To make a support vector machine, we will modify the first term of the cost function $-\log(h_\theta(x)) = -\log\left(\frac{1}{1 + e^{-\theta^T x}}\right)$ so that when $\theta^T x$ (from now on, we shall refer to this as z) is **greater than 1**, it outputs 0. Furthermore, for values of z less than -1, we shall use a straight decreasing line instead of the sigmoid curve. (In the literature, this is called a hinge loss (https://en.wikipedia.org/wiki/Hinge_loss) function.)

Large Margin Intuition

A useful way to think about Support Vector Machines is to think of them as *Large Margin Classifiers*.

If $y=1$, we want $\Theta^T x \geq 1$ (not just ≥ 0)

If $y=0$, we want $\Theta^T x \leq -1$ (not just < 0)

Now when we set our constant C to a very **large** value (e.g. 100,000), our optimizing function will constrain Θ such that the equation A (the summation of the cost of each example) equals 0. We impose the following constraints on Θ :

$$\Theta^T x \geq 1 \text{ if } y=1 \text{ and } \Theta^T x \leq -1 \text{ if } y=0.$$

If C is very large, we must choose Θ parameters such that:

$$\sum_{i=1}^m y^{(i)} \text{cost}_1(\Theta^T x) + (1 - y^{(i)}) \text{cost}_0(\Theta^T x) = 0$$

This reduces our cost function to:

$$J(\theta) = C \cdot 0 + \frac{1}{2} \sum_{j=1}^n \Theta_j^2 \\ = \frac{1}{2} \sum_{j=1}^n \Theta_j^2$$

Recall the decision boundary from logistic regression (the line separating the positive and negative examples). In SVMs, the decision boundary has the special property that it is **as far away as possible** from both the positive and the negative examples.

The distance of the decision boundary to the nearest example is called the **margin**. Since SVMs maximize this margin, it is often called a *Large Margin Classifier*.

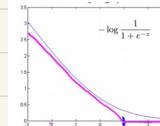
The SVM will separate the negative and positive examples by a **large margin**.

This large margin is only achieved when **C is very large**.

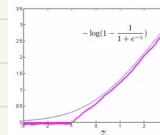
Data is **linearly separable** when a **straight line** can separate the positive and negative examples.

If we have **outlier** examples that we don't want to affect the decision boundary, then we can **reduce C** .

Increasing and decreasing C is similar to respectively decreasing and increasing λ , and can simplify our decision boundary.



Similarly, we modify the second term of the cost function $-\log(1 - h_\theta(x)) = -\log\left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$ so that when z is less than -1, it outputs 0. We also modify it so that for values of z greater than -1, we use a straight increasing line instead of the sigmoid curve.



We shall denote these as $\text{cost}_1(z)$ and $\text{cost}_0(z)$ respectively, note that $\text{cost}_1(z)$ is the cost for classifying when $y=1$, and $\text{cost}_0(z)$ is the cost for classifying when $y=0$, and we may define them as follows (where k is an arbitrary constant defining the magnitude of the slope of the line):

$$z = \theta^T x \\ \text{cost}_1(z) = \max(0, k(1 + z)) \\ \text{cost}_0(z) = \max(0, k(1 - z))$$

Recall the full cost function from (regularized) logistic regression:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m y^{(i)} (-\log(h_\theta(x^{(i)}))) + (1 - y^{(i)}) (-\log(1 - h_\theta(x^{(i)}))) + \frac{\lambda}{2m} \sum_{j=1}^n \Theta_j^2$$

Note that the negative sign has been distributed into the sum in the above equation.

We may transform this into the cost function for support vector machines by substituting $\text{cost}_0(z)$ and $\text{cost}_1(z)$:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{\lambda}{2m} \sum_{j=1}^n \Theta_j^2$$

We can optimize this a bit by multiplying this by m (thus removing the m in the denominators). Note that this does not affect our optimization, since we're simply multiplying our cost function by a positive constant (for example, minimizing $(u - 5)^2 + 1$ gives us 5, multiplying it by 10 to make it $10(u - 5)^2 + 10$ still gives us 5 when minimized).

$$J(\theta) = \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{\lambda}{2} \sum_{j=1}^n \Theta_j^2$$

Furthermore, convention dictates that we regularize using a factor C , instead of λ , like so:

$$J(\theta) = C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=1}^n \Theta_j^2$$

This is equivalent to multiplying the equation by $C = \frac{1}{\lambda}$, and thus results in the same values when optimized. Now, when we wish to regularize more (that is, reduce overfitting), we decrease C , and when we wish to regularize less (that is, reduce underfitting), we increase C .

Finally, note that the hypothesis of the Support Vector Machine is not interpreted as the probability of y being 0 or 1 (as it is for the hypothesis of logistic regression). Instead, it outputs either 1 or 0. (In technical terms, it is a discriminant function.)

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

Mathematics Behind Large Margin Classification (Optional)

Vector Inner Product

Say we have two vectors, u and v :

$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

The length of vector v is denoted $\|v\|$, and it describes the line on a graph from origin $(0,0)$ to (v_1, v_2) .

The length of vector v can be calculated with $\sqrt{v_1^2 + v_2^2}$ by the Pythagorean theorem.

The projection of vector v onto vector u is found by taking a right angle from u to the end of v , creating a right triangle.

- p : length of projection of v onto the vector u .
- $u^T v = p \cdot \|u\|$

Note that $u^T v = \|u\| \cdot \|v\| \cos \theta$ where θ is the angle between u and v . Also, $p = \|v\| \cos \theta$. If you substitute p for $\|v\| \cos \theta$, you get $p = \|v\| - \|u\|$.

So the product $u^T v$ is equal to the length of the projection times the length of vector u .

In our example, since u and v are vectors of the same length, $u^T v = \|u\| \cdot \|v\|$.

$$u^T v = v^T u = p \cdot \|u\| = u_1 v_1 + u_2 v_2$$

If the angle between the lines for v and u is **greater than 90 degrees**, then the projection p will be **negative**.

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \Theta_j^2 \\ = \frac{1}{2} (\Theta_1^2 + \Theta_2^2 + \dots + \Theta_n^2) \\ = \frac{1}{2} \sqrt{(\Theta_1^2 + \Theta_2^2 + \dots + \Theta_n^2)^2} \\ = \frac{1}{2} \|\Theta\|^2$$

We can use the same rules to rewrite $\Theta^T x^{(i)}$:

$$\Theta^T x^{(i)} = p^{(i)} \cdot \|\Theta\| = \Theta_1 x_1^{(i)} + \Theta_2 x_2^{(i)} + \dots + \Theta_n x_n^{(i)}$$

So we now have a new **optimization objective** by substituting $p^{(i)} \cdot \|\Theta\|$ in for $\Theta^T x^{(i)}$:

If $y=1$, we want $p^{(i)} \cdot \|\Theta\| \geq 1$

If $y=0$, we want $p^{(i)} \cdot \|\Theta\| \leq -1$

The reason this causes a "large margin" is because: the vector for Θ is perpendicular to the decision boundary. In order for our optimization objective (above) to hold true, we need the absolute value of our projections $p^{(i)}$ to be as large as possible.

If $\Theta_0 = 0$, then all our decision boundaries will intersect $(0,0)$. If $\Theta_0 \neq 0$, the support vector machine will still find a large margin for the decision boundary.

Kernels I

Kernels allow us to make complex, non-linear classifiers using Support Vector Machines.

Given x , compute new feature depending on proximity to landmarks $l^{(1)}$, $l^{(2)}$, $l^{(3)}$.

To do this, we find the "similarity" of x and some landmark $l^{(i)}$:

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

This "similarity" function is called a **Gaussian Kernel**. It is a specific example of a kernel.

The similarity function can also be written as follows:

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

There are a couple properties of the similarity function:

$$\text{if } x \approx l^{(i)}, \text{ then } f_i = \exp\left(-\frac{\approx 0^2}{2\sigma^2}\right) \approx 1$$

$$\text{if } x \text{ is far from } l^{(i)}, \text{ then } f_i = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$

In other words, if x and the landmark are close, then the similarity will be close to 1, and if x and the landmark are far away from each other, the similarity will be close to 0.

Each landmark gives us the features in our hypothesis:

$$\begin{aligned} l^{(1)} &\rightarrow f_1 \\ l^{(2)} &\rightarrow f_2 \\ l^{(3)} &\rightarrow f_3 \\ &\dots \end{aligned}$$

$$h_\Theta(x) = \Theta_0 f_0 + \Theta_1 f_1 + \Theta_2 f_2 + \Theta_3 f_3 + \dots$$

σ^2 is a parameter of the Gaussian Kernel, and it can be modified to increase or decrease the **drop-off** of our feature f_i . Combined with looking at the values inside Θ , we can choose these landmarks to get the general shape of the decision boundary.

Choosing SVM Parameters

$$\text{Choosing } C \text{ (recall that } C = \frac{1}{\lambda}$$

- If C is large, then we get higher variance/lower bias
- If C is small, then we get lower variance/higher bias

The other parameter we must choose is σ^2 from the Gaussian Kernel function:

With a large σ^2 , the features f_i vary more smoothly, causing higher bias and lower variance.

With a small σ^2 , the features f_i vary less smoothly, causing lower bias and higher variance.

Using An SVM

There are lots of good SVM libraries already written. A. Ng often uses "liblinear" and "libsvm". In practical application, you should use one of these libraries rather than rewrite the functions.

In practical application, the choices you do need to make are:

- Choice of parameter C
- Choice of kernel (similarity function)
- No kernel ("linear" kernel) -- gives standard linear classifier
- Choose when n is large and when m is small
- Gaussian Kernel (above) -- need to choose σ^2
- Choose when n is small and m is large

The library may ask you to provide the kernel function.

Note: do perform feature scaling before using the Gaussian Kernel.

Note: not all similarity functions are valid kernels. They must satisfy "Mercer's Theorem" which guarantees that the SVM package's optimizations run correctly and do not diverge.

You want to train C and the parameters for the kernel function using the training and cross-validation datasets.

Kernels II

One way to get the landmarks is to put them in the **exact same** locations as all the training examples. This gives us m landmarks, with one landmark per training example.

Given example x :

$$f_1 = \text{similarity}(x, l^{(1)}), f_2 = \text{similarity}(x, l^{(2)}), f_3 = \text{similarity}(x, l^{(3)}), \text{and so on.}$$

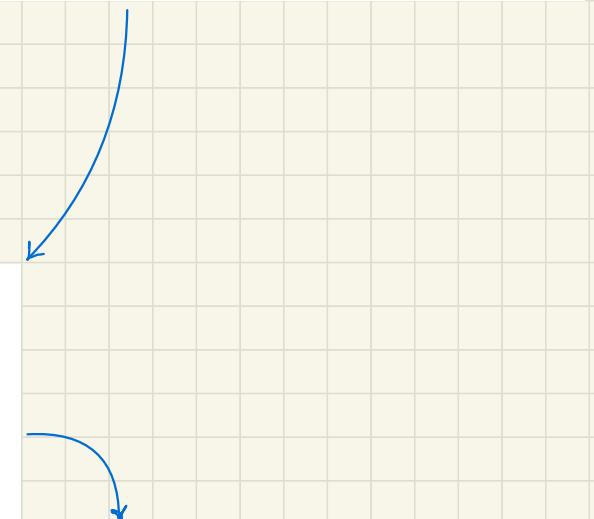
This gives us a "feature vector" $f_{(i)}$ of all our features for example $x_{(i)}$. We may also set $f_0 = 1$ to correspond with Θ_0 . Thus given training example $x_{(i)}$:

$$x^{(i)} \rightarrow [f_1^{(i)} = \text{similarity}(x^{(i)}, l^{(1)}) \ f_2^{(i)} = \text{similarity}(x^{(i)}, l^{(2)}) \ f_3^{(i)} = \text{similarity}(x^{(i)}, l^{(3)})]$$

Now to get the parameters Θ we can use the SVM minimization algorithm but with $f^{(i)}$ substituted in for $x^{(i)}$:

$$\min_{\Theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\Theta^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\Theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^n \Theta_j^2$$

Using kernels to generate $f(i)$ is not exclusive to SVMs and may also be applied to logistic regression. However, because of computational optimizations on SVMs, kernels combined with SVMs is much faster than with other algorithms, so kernels are almost always found combined only with SVMs.



Multi-class Classification

Many SVM libraries have multi-class classification built-in.

You can use the *one-vs-all* method just like we did for logistic regression, where $y \in \{1, 2, 3, \dots, K\}$ with $\Theta^{(1)}, \Theta^{(2)}, \dots, \Theta^{(K)}$. We pick class i with the largest $(\Theta^{(i)})^T x$.

Logistic Regression vs. SVMs

If n is large (relative to m), then use logistic regression, or SVM without a kernel (the "linear kernel")

If n is small and m is intermediate, then use SVM with a Gaussian Kernel

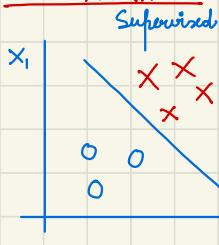
If n is small and m is large, then manually create/add more features, then use logistic regression or SVM without a kernel.

In the first case, we don't have enough examples to need a complicated polynomial hypothesis. In the second example, we have enough examples that we may need a complex non-linear hypothesis. In the last case, we want to increase our features so that logistic regression becomes applicable.

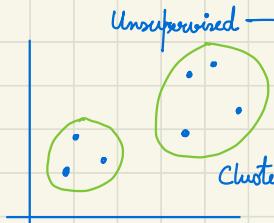
Note: a neural network is likely to work well for any of these situations, but may be slower to train.

UNSUPERVISED LEARNING

1. CLUSTERING



vs



Unsupervised

find me a pattern

Clustering algorithm

market segmentation

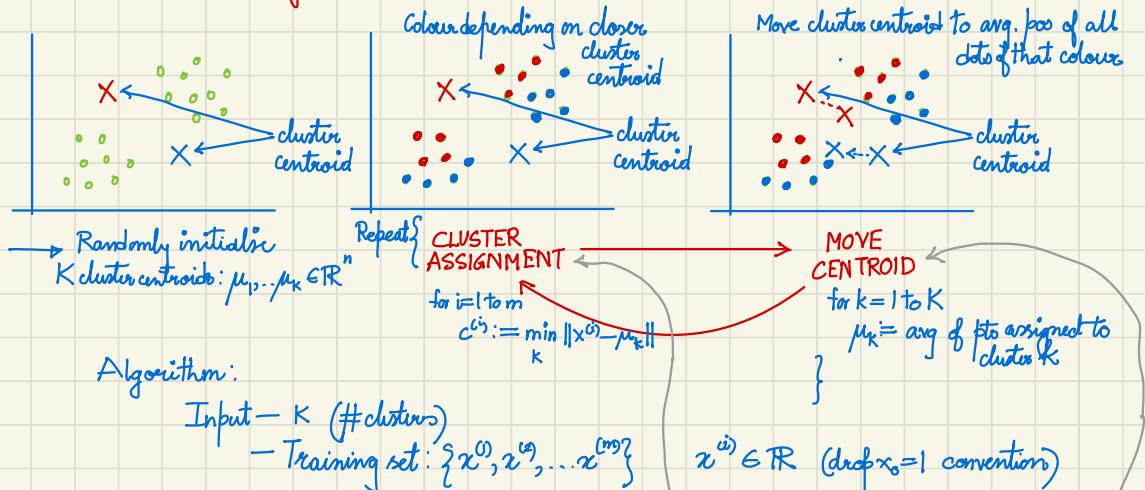
social network analysis

organise computer clusters

galaxy formation

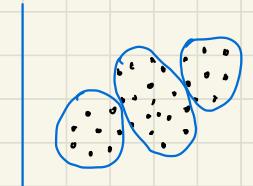
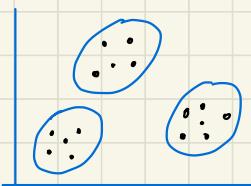
$$\text{Training set: } \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\} \rightarrow \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$$

→ K-Means Algorithm



→ K-Means for non-separated clusters

→ T-shirt sizing



$$\min J(\dots) \text{ wrt } c^{(1)}, \dots, c^{(m)} \\ (\text{holding } \mu_1, \dots, \mu_K \text{ fixed})$$

$$\min J(\dots) \text{ wrt } \mu_1, \dots, \mu_K \\ (\text{holding } c^{(1)}, \dots, c^{(m)} \text{ fixed})$$

OPTIMISATION OBJECTIVE

K-means: $\{c^{(1)}, \dots, c^{(m)}\} \rightarrow \mu_{c^{(1)}} \in \mathbb{R}^n$

loc. of cluster centroid
to which $x^{(i)}$ has been
assigned.

Distortion cost function

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

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

Random Initialisation

- $K < m$
- Pick random K training egs
- set μ_1, \dots, μ_K = these K examples

Avoiding Local optima

$\times 50 - 1000 \rightarrow$ Compute $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$
pick clustering with lowest J

Choosing K:



→ Benchmark by downstream performance
eg: T-shirt size business

2. DIMENSIONALITY REDUCTION

MOTIVATIONS

→ **Data Compression**: Reduce data from 2D \rightarrow 1D

$$\begin{aligned} x^{(1)} \in \mathbb{R}^2 &\xrightarrow{\quad} z^{(1)} \in \mathbb{R} \\ x^{(m)} \in \mathbb{R}^2 &\xrightarrow{\quad} z^{(m)} \in \mathbb{R} \end{aligned} \quad \text{project onto reduced dimension}$$

→ reduces disk space/storage

→ faster algorithm

$$k = \# \text{PC}$$



→ **Data Visualisations**:

Dimensional reduction: $\{x^{(1)}, \dots, x^{(m)}\} \rightarrow \{z^{(1)}, \dots, z^{(m)}\}$

$$z^{(1)} \in \mathbb{R}^n$$

$$z^{(1)} \in \mathbb{R}^k$$

$$k \leq n$$

→ $k=2, 3$ (plot to visualize correlations)

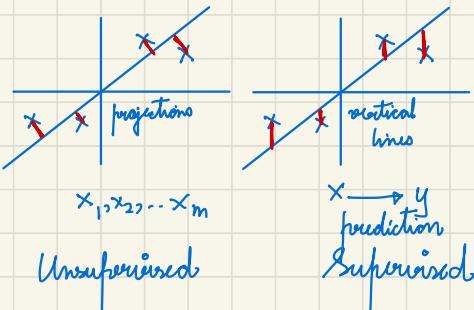
→ **Analysis algorithm**: PCA

→ Problem formulation: find lower-D surface to project data onto by minimising squared projected error.

$\rightarrow n\text{-D} \rightarrow k\text{-D}$ $[u^{(1)}, \dots, u^{(k)}]$: min Proj Error

↓
project onto linear subspace of these vectors

PCA \neq Linear Regression



→ Data preprocessing

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

Preprocessing: feature scaling / mean normalisation

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$x_j^{(i)} \rightarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

$$U_{\text{reduce}} = z = \begin{bmatrix} | & | \\ u^{(1)} & \dots & u^{(k)} \\ | & | \end{bmatrix}^T x$$

$$= \begin{bmatrix} | & | \\ (u^{(1)})^T & \dots & (u^{(k)})^T \\ | & | \end{bmatrix} x$$

→ Algorithm: $n\text{-D} \rightarrow k\text{-D}$

- Compute "covariance matrix" $\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)}) (x^{(i)})^T$ $n \times n$
- Compute "eigenvectors" of Σ

$$[U, S, V] = \text{svd}(\Sigma)$$

or find library for corresponding language

$$U = \begin{bmatrix} | & | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(m)} \\ | & | & | & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

K $n \rightarrow k$

Summary: → Mean normalisation + Feature Scaling
 $\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)}) (x^{(i)})^T$

PCA

Algorithm

$$[U, S, V] = \text{svd}(\Sigma)$$

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

$$x \in \mathbb{R}^n$$

$$z = U_{\text{reduce}} * x$$

Compression: $z = U_{\text{reduce}}^T x$ $n-D \rightarrow k-D$

Decompression: $k-D \rightarrow n-D$

$$x_{\text{approx}} = U_{\text{reduce}} \cdot z$$

Choosing k :

Avg sq. Proj. Error: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$

Total var. in data: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$

Choose k : smallest value s.t.

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

99% is retained

$$\frac{1}{m} \sum \|x^{(i)}\|^2$$

Algorithm:

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

$$S = \begin{bmatrix} S_{11} & & \\ & S_{22} & \\ & & \ddots \\ & & & S_{nn} \end{bmatrix}$$

For given k : Check if

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

→ Increase k

Extract inputs:

Unlabeled training set: $x^{(1)}, \dots, x^{(m)} \in \mathbb{R}^{10000}$ (say)

$$z^{(1)}, \dots, z^{(m)} \in \mathbb{R}^{1000}$$

↓ PCA on training set → mapping can be applied as well on
 $x_{cv}^{(i)}$ or $x_{test}^{(i)}$

→ New training set:

$$\{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\}$$

APPLICATION OF PCA

- Compression
 - Reduce memory/disk needed to store data
 - Speed up learning algorithm
- Visualisation]
k = 2 or 3

~~BAD USE OF PCA: Prevent Overfitting: $x^{(i)} \rightarrow z^{(i)}$ $k < n \rightarrow$ less likely to overfit~~

→ might work OK. Instead use REGULARISATION if $g_{\theta}^T x$ is retained it might throw away imp data

Design of ML System: Before implementing PCA, consider using original data $x^{(i)}$. If that doesn't work, consider dimensional reduction.

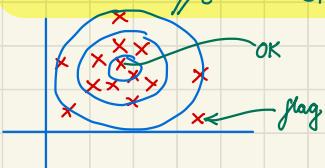
ANOMALY DETECTION

Problem: Given dataset $\{x^{(1)}, \dots, x^{(m)}\}$

Is x_{test} anomalous?

→ Build model $p(x)$

If $p(x_{\text{test}}) < \varepsilon \rightarrow$ Flag anomaly
 $> \varepsilon \rightarrow$ OK



Example: FRAUD DETECTION

- $x^{(i)}$ = features of user i's activities
- Model $p(x)$ from data
- Id unusual users by checking $p(x) < \varepsilon$

MANUFACTURING

MONITORING COMPUTERS IN DATA CENTRES

- $x^{(i)}$: features of machine i
- $p(x) < \varepsilon$ to predict machines that might go down

ANOMALY DETECTION ALGORITHM WITH GAUSSIAN DIST

Training set: $\{x^{(1)}, \dots, x^{(m)}\}; x \in \mathbb{R}^n$

Assume: $p(x) = p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) \dots p(x_m; \mu_m, \sigma_m^2)$ [IID]

Algorithm:

- Choose features x_i which might be indicative of anomalous examples

• Fit parameters: $\mu_1, \dots, \mu_m, \sigma_1^2, \dots, \sigma_m^2$

$$x^{(i)} \sim p(x^{(i)}; \mu_i, \sigma_i^2) \quad \mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

- Given new example x , compute $p(x)$: $p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$



→ Developing/Evaluating an Anomaly Detection System

→ w/ real # evaluation $y=1$ $y=0$

→ assume labelled data of anomalous/non-anomalous egs

→ training set: $\{x^{(1)}, \dots, x^{(m)}\}$ assume non-anomalous

→ define CV and test sets

Motivating example: Aircraft engines



Algorithm:

→ Fit $p(x)$ on training set: $\{x^{(1)}, \dots, x^{(m)}\}$

→ On CV/test example, x , 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 ϵ

Classification accuracy not great ← skewed data set

instead use

Possible Eval. Metrics

True +ve, False +ve
True -ve, False -ve

Precision/Recall

F1 Score

ANOMALY DETECTION vs SUPERVISED LEARNING

→ very few $y=1$, +ve egs
very large, $y=0$, -ve egs

→ large # of +ve or -ve egs

→ many different types of anom.
↳ hard for any algorithm
to learn from +ve egs

→ enough +ve egs to get
a sense of what +ve egs
are like → generalizable
to future

→ future anom. might look
nothing like anom. egs
seen so far

Very few $y=1$

EXAMPLES

$y=0, y=1$, almost same

- Fraud detection
- Manufacturing
- monitoring machines
in data clusters

- Email spam classification
- Weather prediction
- cancer classification

Choosing Features:

→ Non-Gaussian features →

func "transf" n To Gauss. func^{o.05}

X → log(X) or X^{0.05}

→ Error analysis

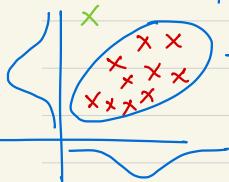
Wish: $p(x)$ large for normal egs
small for anom egs

Reality: comparable

For anomaly not being flagged, get
inspiration for a new feature

→ choose features which will take large
or small values

ANOMALY DETECTION w/ MULTIVARIATE GAUSSIAN



→ doesn't assume covariance is diagonal with equal elements

→ multivariate Gaussian distⁿ

$$\sigma^2 \rightarrow \Sigma$$

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

→ apply to anomaly detection algorithm

Parameters: μ, Σ

Parameter fitting: Given training set,
 $\{x^{(1)}, \dots, x^{(m)}\} \quad x \in \mathbb{R}^n$

$$\text{Estimate: } \mu = \frac{1}{m} \sum x^{(i)}$$

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

→ Fit model $p(x)$ by setting $\mu = \dots, \Sigma = \dots$

→ Given new example x , compute $p(x)$

→ Flag anomaly if $p(x) < \epsilon$

→ Original model

$$p(x_1; \mu_1, \sigma_1^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$$

Manually create features to capture anomalies where x_1, x_2 take unusual combinations of values.

$$x_3 = \frac{x_1}{x_2} = \frac{\text{CPU load}}{\text{memory}}$$

→ Computationally cheaper (alternatively, scales better to large n)
 $n=10,000, m=100,000$

OK even if m (training set size) is small

vs. → Multivariate Gaussian

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

→ Automatically captures correlations between features

$$\Sigma \in \mathbb{R}^{n \times n}$$

$$\Sigma^{-1}$$

Computationally more expensive

$$\Sigma$$

$$n \times n$$

$$\rightarrow x_1 = x_2 = x_3 = x_4 = x_5 = \dots$$

Must have $m > n$ or else Σ is non-invertible.

Andrew Ng

RECOMMENDER SYSTEMS

Example: Predicting movie ratings

n_u : #users

n_m : #movies

$r_{(i,j)} = 1$ if user j has rated movie i

$y_{(i,j)} = \text{rating given}$

ALGORITHM

→ Content-based recommender system

$$x_1 \rightarrow \text{romance} \in \{0, 1\} + x_0 = 1$$

$$x_2 \rightarrow \text{action} \in \{0, 1\}$$

For each user j , learn parameters $\theta^{(j)} \in \mathbb{R}^{n+1}$

Predict user j 's rating for movie i w/ $(\theta^{(j)})^T x^{(i)}$ stars

LINEAR REGRESSION

for each user

Problems: given from data already rated what would a user have rated a movie not seen?

Problem Formulation: $r(i,j), y^{(i,j)}$
 $\theta^{(j)}, x^{(i)}$

$m^{(j)} = \# \text{movies rated by user } j$

To learn $\theta^{(j)}$: $\min_{\theta^{(j)}} \frac{1}{m^{(j)}} \sum_{i: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^{n_u} (\theta_k^{(j)})^2$ [optimization formulation]

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

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

$J(\theta)$

Gradient Descent Update:

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

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left[\dots + \lambda \theta_k^{(j)} \right] \quad k \neq 0$$

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

Basically Linear Regression w/o $\frac{1}{m}$ term

→ Collaborative Filtering:

learn features on its own

↳ optimisation algorithm

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

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

Collaborative filtering

[Given $\theta \rightarrow x$]

[Given $x \rightarrow \theta$]

[Given $\theta \rightarrow x \rightarrow \theta \rightarrow x \dots$]

Collaborating filtering

→ Collaborative filtering algorithm

Solve simultaneously

Collaborative filtering algorithm

→ 1. Initialize $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ to small random values.

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

$$x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j: r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right) \quad \cancel{x \in \mathbb{R}^n, \theta \in \mathbb{R}^m}$$

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

Vectorised:

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

Predicted Ratings: $x \odot \theta^T$

Low Rank Matrix Factorisation

→ Finding Related Movies:

For pdt i , learn feature vector $x^{(i)} \in \mathbb{R}^n$

Relation b/w movies $i \rightarrow j$

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

3. For a user with parameters θ and a movie with (learned) features x , predict a star rating of $\theta^T x$.

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

+ MEAN NORMALISATION

LEARNING WITH LARGE DATASETS

$m \sim O(10^6 - 10^7)$ (typical) \downarrow

(Batch) Gradient descent: $\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m \text{sum over } O(10^6 - 10^7) \text{ sample fr each iteration} !!$

use all m examples in each iterations

Perhaps select subset $\sim O(10^3)$

Sanity check

plot Learning curve

c_v

c_y

train

high var \times

Increasing m

high bias \checkmark

→ Stochastic Gradient descent: → use 1 eg in each iteration

1. Randomly shuffle/reorder training examples

2. Repeat

for $i = 1, \dots, m$

$$\theta_j := \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \forall j = 1, \dots, n$$

}

Checking for convergence

compute cost $(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$ before updating θ using $x^{(i)}$

every 1000 iter, plot cost ← averaged over 1000 cgs

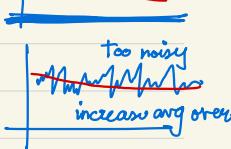
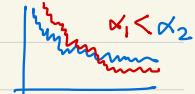
→ Mini-batch gradient descent → use b examples in each iteration

↓
mini batch size

stochastic gradient descent ← iff vectorization

for $i = 1, 11, \dots, 991$

if $b = 10, m = 1000$



Adapting Stochastic GD for ONLINE LEARNING

- unlimited data
- can adapt to changing user preferences
- 1 eg at a time

The PHOTO OCR Problem → computers to read text in images we take

Given

Photo,

- Find where texts appear
- Id/Read them

Optical Character Recognition

Uses

- Search photos using text
- camera for blind
- car navigation

PHOTO OCR PIPELINE

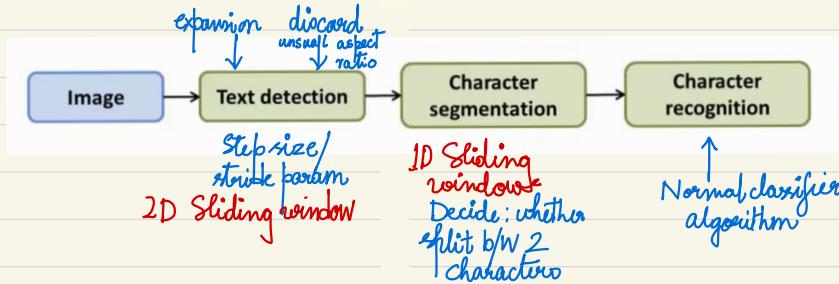
MACHINE LEARNING
PIPELINE

- text detection
- character segmentation
- character classification

A system with many stages/components
several of which may use machine learning.

→ SLIDING WINDOWS CLASSIFIER

Pedestrian detection → Text detection



→ Synthetic data: Font library + random background

↗ unlimited supply of labelled data

→ Introducing distortions * CAVEAT: Distortions REPRESENTATIVE of egs in test set
→ not meaningless distortion

Speech Recognition

→ original audio + include distorted background
eg, bad cellphone
traffic background

1 training eg
many training egs

Discussion on getting more data

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

+ CEILING ANALYSIS

- Giving right answers
 ↳ has high variance/bias
- Id stage showing max. improvement
 ↓
 work more on this stage

Summary: Main topics

Supervised Learning

- Linear regression, logistic regression, neural networks, SVMs

$(x^{(i)}, y^{(i)}) \rightarrow \text{Labelled data}$

Unsupervised Learning

- K-means, PCA, Anomaly detection

$x^{(i)} \rightarrow \text{unlabelled data}$

Special applications/special topics

- Recommender systems, large scale machine learning. (including parallelised)

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.

→ precision/recall

→ F1 score

→ training/cv/test set

DEBUGGING