

Coursera, Stanford

Machine

Learning

Andrew N.G.

Course Summary : Main Topics

- Supervised Learning
 - Linear Regression
 - Logistic Regression
 - Neural Networks
 - Support Vector Machines (SVM)

- Unsupervised Learning

- K-means
 - PCA
 - Anomaly Detection

- Special Applications / special topics

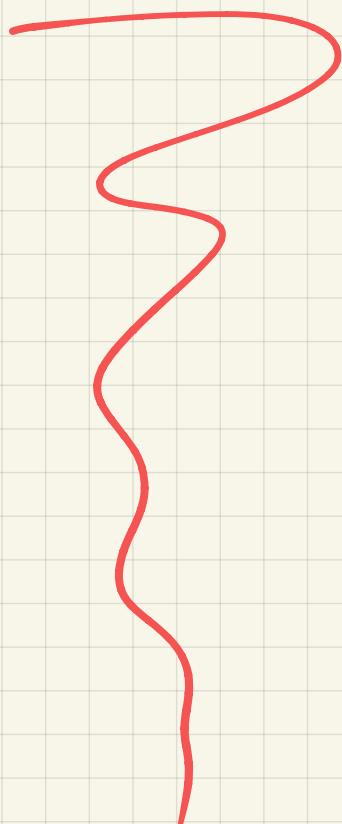
- Recommender Systems
 - Large Scale machine learning

- Advice on building a ML System

- Bias / Variance
 - Regularization
 - Deciding what to work on next:
 - Evaluation of Learning Algorithms
 - Learning curves
 - Error Analysis
 - Ceiling Analysis

Disclaimer: All the notes written below were strongly based on the course material.

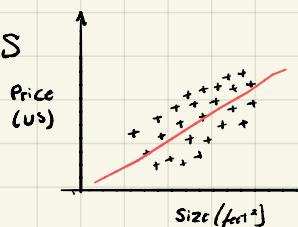
Airton Tiago → 22-04-2020



Week 1 - Linear Regression with one Variable

Model Representation

- Housing Prices



Training set of housing prices
Notation

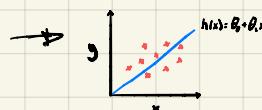
m : Number of training examples

$x^{(i)}$: "Input" variable / features

$y^{(i)}$: "Output" variable / "Target" variable

$x^{(1)}$: 1st line of data set (training)

- How do we represent h ? $h_{\theta}(x) = \theta_0 + \theta_1 x$



Supervised Learning → Given the "right answer" for each example in the data

Regression Problem → Predict real-valued output

Classification Problem → Predict discrete-valued output

Training Set

Learning Algorithm

x
Size of house
 $\rightarrow h$ (hypothesis)
 y
Estimated price

h maps from x 's to y 's

Linear Regression with one Variable

Univariate Linear Regression

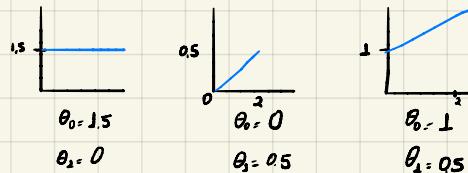
fancy way to say "one variable"

COST function

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

θ_0, θ_1 's = Parameters → How to choose them?

(Choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training examples $(x^{(i)}, y^{(i)})$)

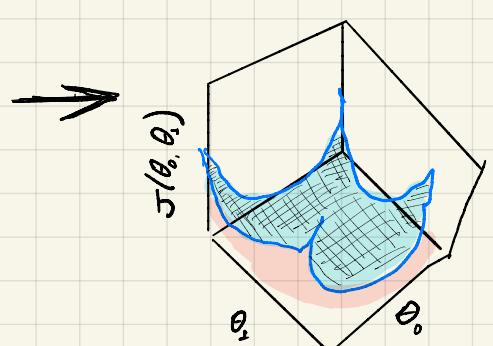
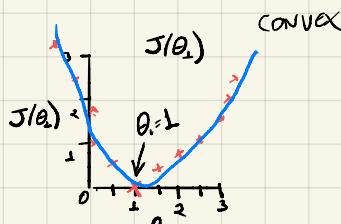
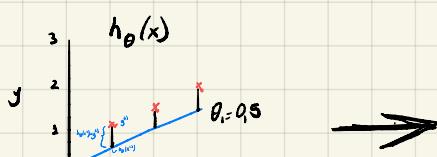


So, we need to solve a minimization problem: $\underset{\theta_0, \theta_1}{\text{minimize}} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \rightarrow \underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$$

cost function
Squared error function

Cost function Intuition



For fixed θ_1 , this is a function of θ_0

$$J(0.5) = \frac{1}{2m} [(0.5-1.0)^2 + (1-2)^2 + (1.5-3)^2] = \frac{1}{2m} [3.5] \approx 0.58$$

function of the parameter θ_1

Gradient Descent

We will use this Gradient to descend to lowest point in the Cost function. When you have two or more derivatives of the same function, they are called Gradient.

Repeat until convergence $\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \text{ (for } j=0 \text{ and } j=1\text{)} \}$

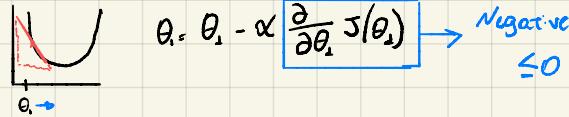
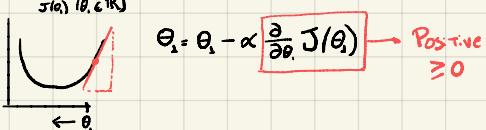
Correct: Simultaneous update

Tangent derivative (tangential line to a function) → determining direction

The slope of the tangent is the derivative at that point and it will give us a direction to move towards.

$$\begin{aligned} \text{Temp 0} &:= \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \text{Temp 1} &:= \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \end{aligned} \rightarrow \begin{aligned} \theta_0 &:= \text{Temp 0} \\ \theta_1 &:= \text{Temp 1} \end{aligned} \text{ := Assignment}$$

$J(\theta)$ (U-shaped)



Learning Rate α

- If α is too small, gradient descent can be SLOW
- If α is too large, gradient descent can overshoot the minimum. It may fail to converge or even diverge
- If θ_1 is initialized in the local minimum the slope is equal to zero $\rightarrow \theta_1 := \theta_1 - \alpha \left(\frac{\partial}{\partial \theta_1} J(\theta_1) \right) = 0 \rightarrow \theta_1 := \theta_1$

Gradient descent can converge to a local minimum, even with the learning rate α fixed

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.

Gradient descent Algorithm

Repeat until convergence $\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \end{array} \right\}$

update θ_0 and θ_1

SIMULTANEOUSLY

Convex function = Bowl-shaped

There is no local minima, just global optimum \rightarrow for Linear Regression, in specific choice of cost function $J(\theta_0, \theta_1)$.

"Batch" Gradient Descent

Sometimes is called by:

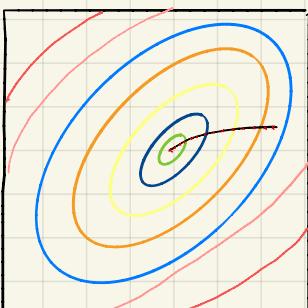
\hookrightarrow "Batch": each step of gradient descent uses all the training examples

\hookrightarrow All the time were looking to this: $\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$

\hookrightarrow sum all over for each batch

- Gradient descent performs better than normal equation when using LARGE datasets.

While gradient descent can be susceptible to local minima in general, the optimization problem we've proposed above for linear regression has only one global, and no other local, optima: thus gradient descent always converges (assuming the learning rate is not too large) to the global minimum. Indeed, J is a convex quadratic function.



The ellipses shown in the left are the contours of a quadratic function. Also shown is the trajectory taken by gradient descent, which was initialized at (10, 30). The x 's in the figure (joined by straight lines) mark the successive values of θ that gradient descent went through as it converged to its minimum.

Linear Algebra Review

• Matrices and Vectors

Definition(s):

- Matrix = Rectangular array of numbers
- Dimension = Number of rows \times Number of columns
- Matrix Elements = Entries of matrix.
 A_{ij} = "i, j entry" in the ith row, jth column
- Vector = An m_x matrix $\rightarrow y \cdot \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

• Matrix Multiplication

$$A = \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix} \rightarrow A_{2 \times 2}$$

$$B = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \rightarrow B_{2 \times 3}$$

$$A_{2 \times 2} \times B_{2 \times 3} = C_{2 \times 3}$$

has to be the same
dimension of entries

$$\rightarrow \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = \begin{bmatrix} (2 \times 1 + 1 \times 4) & (2 \times 2 + 1 \times 5) & (2 \times 3 + 1 \times 6) \\ (1 \times 1 + 3 \times 4) & (1 \times 2 + 3 \times 5) & (1 \times 3 + 3 \times 6) \end{bmatrix} = \begin{bmatrix} 6 & 9 & 12 \\ 13 & 17 & 21 \end{bmatrix}$$

• Properties

$$\rightarrow A \times B \neq B \times A \text{ (Not commutative)}$$

$$\hookrightarrow \text{Identity Matrix} = \text{denoted } I \text{ (or } I_{n \times n})$$

$$\hookrightarrow \text{Examples: } \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \rightarrow \text{for any matrix } A, A \cdot I = I \cdot A = A$$

• Inverse and Transpose

- Not all numbers have an inverse

Matrix Inverse

\hookrightarrow If A is an m_xn matrix, and it has an inverse, $A(A^{-1}) = A^{-1}A = I$

Matrices that don't have an inverse are "singular" or "degenerate"

• Matrix Transpose

$$A = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 5 & 9 \end{bmatrix} \quad B = A^T = \begin{bmatrix} 1 & 3 \\ 2 & 5 \\ 0 & 9 \end{bmatrix}$$

Let A be an m_xn matrix, and let $B = A^T$
Then B is an n_xm matrix, and $B_{ij} = A_{ji}$

Week 2 - Linear Regression with Multiple Variables

- Multiple features

Previously: $h_{\theta}(x) = \theta_0 + \theta_1 x_1$ $\xrightarrow{\text{Now}}$ $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \rightarrow$ for convenience of notation, define $x_0 = 1$ ($x_0^{(i)} = 1$)

Gradient Descent for multiple variables

Hypothesis: $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Parameters: $\theta_0, \theta_1, \dots, \theta_n$

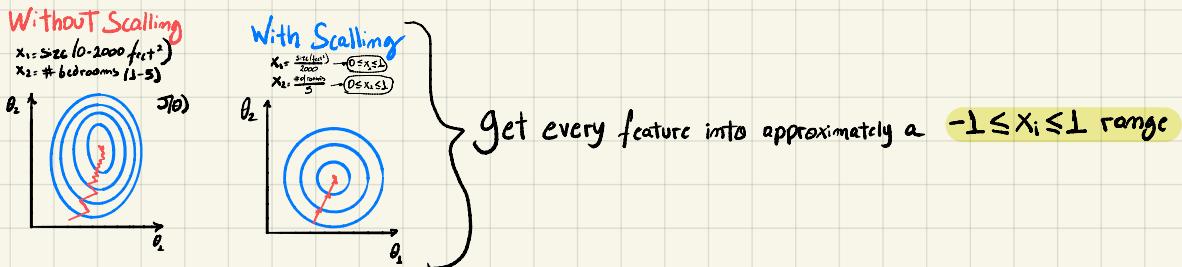
Cost function: $J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Gradient Descent: Repeat: $\left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \right\}$ (Simultaneously update for every $j = 0, \dots, n$)

Now we have $m=2$, before it was $m=1$

Gradient Descent in practice

feature Scaling: Make sure features are on similar scale



Mean Normalization = Replace x_i with $x_i - \mu_i$ to make features have approximately zero mean (Do not apply to $x_0 = 1$)

$$\xrightarrow{\text{Average value of } x \text{ in training set}} x_1 \leftarrow \frac{x_1 - \mu_1}{s_1} \quad x_2 \leftarrow \frac{x_2 - \mu_2}{s_2}$$

(s_1 range/min-max or standard deviation)

How to make sure gradient descent is working correctly?

→ "Debugging":

Example automatic convergence test:

→ Declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration

→ $J(\theta)$ should decrease after every iteration.

→ for sufficiently small α , $J(\theta)$ should decrease on every iteration

→ But if α is too small, gradient descent can be slow to converge.

· If α is too small: slow convergence.

· If α is too large: $J(\theta)$ may not decrease on every iteration; may not converge (slow convergence also possible).

features and polynomial regression

$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$, in the house price prediction example x_1, x_2, x_3 would be

$$\begin{cases} x_1 = (\text{size}) \\ x_2 = (\text{size})^2 \\ x_3 = (\text{size})^3 \end{cases}$$

Normal Equation

Method to solve for θ analytically

Singular = Non-invertible matrix

Examples: $m = 4$.

x_0	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$ $m \times (n+1)$

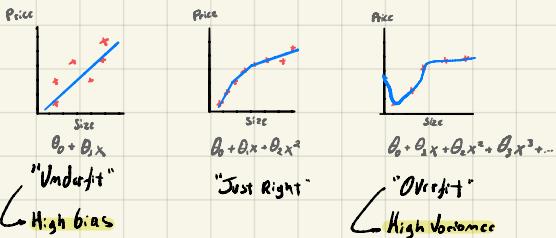
$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$ m -dimensional vector

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

Andrew T

The problem of Overfitting

Example: Linear Regression (housing Prices)



Overfitting: if we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples.

How well an hypothesis apply to a new example

Addressing Overfitting

1. Reduce number of features

- Manually select which features to keep
- Model selection algorithm

2. Regularization

- Keep all the features, but reduce magnitude/values of parameters θ_j
- Works well when we have a lot of features, each of which contributes a bit to predicting y .
- $J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$

Regularization term
Regularization parameter

If set to large number, for example 10¹⁰, it can lead to "underfit" where $h_\theta(x) = \theta_0$

...
...

If we have overfitting from our hypothesis function, we can reduce the weight that some of the terms in our function carry by increasing their cost

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

We'll want to eliminate the influence of $\theta_3 x^3$ and $\theta_4 x^4$. Without actually getting rid of these features or changing the form of our hypothesis,

We can instead modify our cost function:

$$\text{min}_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

We have added two extra terms at the end to INFLATE the cost of θ_3 and θ_4 .

Week 3 - Logistic Regression

Logistic Regression

- Hypothesis: $h_{\theta}(x) = g(\theta^T x)$

where g is the sigmoid function $g(z) = \frac{1}{1+e^{-z}}$

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

- Cost function

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

- Gradient

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

- Regularized Logistic function

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

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \quad \text{for } j \geq 1$$

Week 4 - Neural Networks: Representation

- Non-Linear hypothesis

- Neurons and the Brain

Neural Networks:

Origins: Algorithms that try to mimic the brain

Was very widely used in 80s and early 90s; popularity diminished in late 90s.

Recent resurgence: State-of-the-art technique for many applications

- Model Representation

Week 5 - Neural Networks: Learning

- Cost function

- Backpropagation Algorithm = Algorithm that try to minimize cost function

$\delta_j^{(l)} = \text{"error" of node } j \text{ in layer } l$.

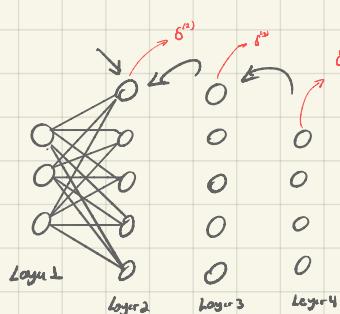
we do not compute it for input layer

for each output unit / layer $L-1$

$$\delta_j^{(L-1)} = a_j^{(L-1)} - y_j$$

$$\delta_j^{(L)} = (\Theta^{(L)})^T \delta_j^{(L-1)} * g'(z^{(L)})$$

$$\delta^{(L)} = (\Theta^{(L)})^T \delta^{(L-1)} * g'(z^{(L)})$$



"Backpropagation" is neural-network terminology for minimizing our cost function

↳ the goal is to compute:

$$\min_{\Theta} J(\Theta)$$

That is we want to minimize our cost function J using an optimal set of parameters in there.

So, in this part of the course we looked at the equations we use to compute the partial derivative of $J(\Theta)$:

$$\frac{\partial}{\partial \Theta^{(l)}_{ij}} J(\Theta)$$

• Backpropagation Algorithm

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

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j)

For $i=1$ to m :

$$\text{Set } a^{(2)} = x^{(i)}$$

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

Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l-1)} + a_j^{(l-1)} \delta_i^{(l)}$

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

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

$$\frac{\partial}{\partial \Theta^{(l)}_{ij}} J(\Theta) = D_{ij}^{(l)}$$

Gradient checking → computed **before** compute backpropagation.

→ will assure that our backpropagation works as intended.

We can approximate the derivative of our cost function with:

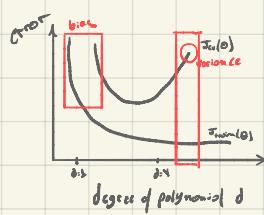
$$\rightarrow \frac{\partial}{\partial \Theta} J(\Theta) \approx \frac{J(\Theta + \epsilon) - J(\Theta - \epsilon)}{2\epsilon} \quad \epsilon \text{ should be } \approx 10^{-4}$$

Week 6 - Advice for Applying ML

- Machine Learning diagnostic
- Evaluate a hypothesis
- Model selection and train/validation/test sets
- Diagnosing bias vs Variance

Suppose your learning algorithm is performing less well than you were hoping. ($J_{cv}(\theta)$ or $J_{test}(\theta)$ is high).

Is it a bias problem or a variance problem?

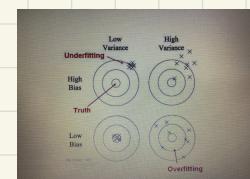
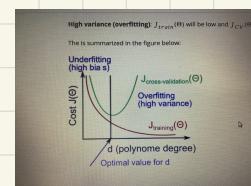


Bias (underfit):

$J_{train} \text{ high}$
 $J_{cv} \text{ high}$
 $J_{train} \approx J_{cv}$

Variance (overfit)

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



- Regularization and bias/variance

• Large λ = high bias/underfit $\rightarrow \lambda = 1000, \theta_0 \approx 0, \theta_1 \approx 0 \rightarrow h(x) \approx 0$

• Intermediate λ =

• Small λ = high Variance/overfit $\lambda = 0 \rightarrow$ nothing changes

for choose the regularization parameter λ properly you can try a range (0, 0.01, 0.02 ... 10) and store each θ^* to see the CV error and choose the best

Bias/variance as a function of the regularization parameter

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

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

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

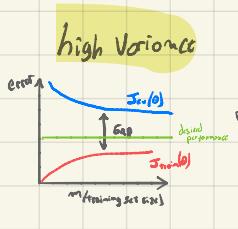
In order to get the "just right" value for λ , we should follow these steps:

1. Create a list of lambdas
2. Create a set of models with different degrees or any other variants.
3. Iterate through the λ s and for each λ go through all the models to learn some Θ
4. Compute the cross validation error using the learned Θ (computed with λ) on the $J_{\text{cv}}(\theta)$ **Without regularization or $\lambda=0$.**
5. Select the best combo that produces the lowest error on the cross validation set
6. Using the best combo Θ and λ , apply it on $J_{\text{train}}(\theta)$ to see if it has a good generalization of the problem.

Learning Curves



If a learning algorithm is suffering from high bias, getting more training data will not help much.



If a learning algorithm is suffering from high variance, getting more training data is likely to help.

Deciding What to Try Next (revisited)

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

Week 6 - ML System Design

Prioritizing what to work on: Spam classification Example

Error Analysis

- Manually pick some wrong predictions in the CV data and try to figure out what happened.

Why is the recommended approach to perform error analysis using the cross validation data and to compute $J_{\text{cv}}(\theta)$ rather than the test data used to compute $J_{\text{test}}(\theta)$?

- If we develop new features by examining the test set, then we may end up choosing features that work well specifically for the test set, so $J_{\text{test}}(\theta)$ is no longer a good estimate of how well we generalize to new examples.

The recommended approach to solving machine learning problems is to:

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

Handling Skewed Data

- Error metrics for **skewed classes**

→ It occurs when we have a lot more examples in one class than another

Precision/Recall

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

Cancer Example:

Precision: Of all predicted as 1, what are the real 1?

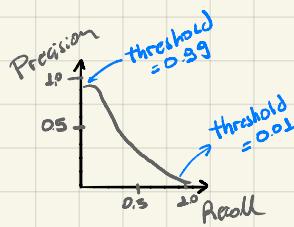
$$\text{Precision} = \frac{\text{True Positive}}{\# \text{predicted positive}} = \frac{\text{True Positives}}{\text{True Positives} + \text{false positives}}$$

Recall: Of all patients that HAVE cancer, what fraction did we correctly detect as having cancer?

$$\text{Recall} = \frac{\text{True positives}}{\# \text{actual positives}} = \frac{\text{True Positives}}{\text{True Positives} + \text{False negatives}}$$

Trading Off Precision and Recall

- for logistic regression: if we want to predict $y=1$ (cancer) only if very confident → **Higher precision Lower recall**
- if we want to avoid missing too many cases of cancer (Avoid False negatives) → **Higher recall Lower precision**



More generally: Predict 1 if $h_\theta(x) \geq \text{threshold}$

Which algorithm is better?

- How to compare precision/recall numbers?

F_1 Score (F Score)

$$F_1 \text{ Score} = 2 \frac{P \cdot R}{P + R}$$

	Precision (P)	Recall (R)	Average	F_1 Score
Algorithm 1	0.5	0.4	0.45	0.444
Algorithm 2	0.7	0.1	0.1	0.175
Algorithm 3	0.02	1.0	0.51	0.0392

→ Predict 1,0 all the time

- A reasonable way to pick the value to use for the threshold in a logistic regression classifier: Measure precision and recall on the **Cross validation set** and choose the value of threshold which maximizes $2 \frac{P \cdot R}{P + R}$

Data for Machine Learning

article [Banko and Brill, 2001] → Classify between confusable words.

"It's not who has the best algorithm that wins. It's who has the most data."

Large data rationale

- Use a learning algorithm with many parameters (e.g. logistic regression/linear regression with many features; neural network with many hidden units) → Low bias → $J_{\text{train}}(\theta)$ will be small
- Use a very large training set (unlike to overfit) → Low variance
 - $J_{\text{train}}(\theta) \approx J_{\text{val}}(\theta)$
 - $J_{\text{val}}(\theta)$ will be small

Week 7 - Support Vector Machines

SVMs are considered by many to be the most powerful 'black box' learning algorithm, and by posing a cleverly-chosen optimization objective, one of the most widely used learning algorithms

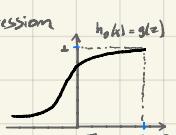
Large Margin Classification

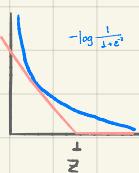
- Optimization Objective

Alternative view of logistic regression

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

If $y=1$, we want $h_\theta(x) \approx 1, \theta^T x \gg 0$
If $y=0$, we want $h_\theta(x) \approx 0, \theta^T x << 0$

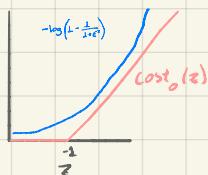




Cost Example: $-(y \log h_0(x) + (1-y) \log(1-h_0(x)))$

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

If we want $\theta^T x > 0$ if we want $\theta^T x < 0$
 $y=1$ $y=0$



Support Vector Machine

Cost function of logistic Regression

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \left(-\log h_0(\theta^T x^{(i)}) \right) + (1-y^{(i)}) \left(-\log(1-h_0(\theta^T x^{(i)})) \right) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

$\text{cost}_1(\theta^T x^{(i)}) \quad \text{cost}_0(\theta^T x^{(i)}) \quad \theta$

Cost function of SVM:

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

C (large constant) $A + \lambda B \rightarrow C = \frac{1}{\lambda} \rightarrow CA + B$

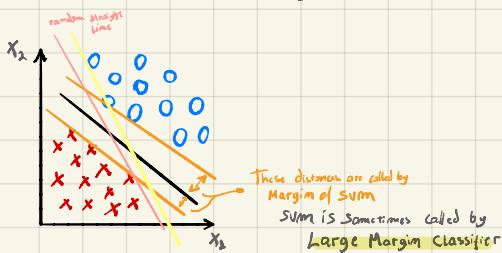
$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

hypothesis: $h_0(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$

They give the same optimal value of θ when $C = \frac{1}{\lambda}$

Large Margin Intuition

SVM Decision Boundary: Linearly Separable Case



SVM try to separate the classes with the large margin possible.

The mathematics behind large margin classification (optional)

Kernels 1

Similarity functions

SVM with Kernels

hypothesis: Given x , computes features $f \in \mathbb{R}^{mn}$

Predict " $y=1$ " if $\theta^T f \geq 0$ $\theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$

Training:

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

SVM PARAMETERS

$C = \frac{1}{\lambda}$ → Large C = Lower bias, high variance → (Small λ) → prone to overfitting
 Small C = Higher bias, low variance → (Large λ) → prone to underfitting

σ^2 → Large σ^2 = features f_i vary more smoothly
 Higher bias, lower variance

Small σ^2 = features f_i vary less smoothly
 Lower bias, higher variance

if using LR



→ higher slope

Suppose you train an SVM and find it overfits your training data.

The reasonable next step is: Decrease C and Increase σ^2

Using AN SVM

Use SVM software packages to solve for parameters θ

Need to specify:

Choice of parameters C → Informally, the C parameter is a positive value that controls the penalty for misclassified training examples. A large C parameter tells SVM to try to classify all the examples correctly. C plays a role similar to $\frac{1}{\lambda}$, where λ is the regularization parameter that we were using previously for logistic regression.

Choice of Kernel (similarity function)

- E.g. No Kernel is also called "Linear Kernel" → Standard linear classifier → # features large, # training samples small
- Predict " $y=1$ " if $\theta^T x \geq 0$

Gaussian Kernel

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

Need to choose σ^2 (sigma squared)

When to use: m is small and/or m is large → $m = \# \text{ features}$
 $m = \# \text{ training examples}$

Kernel (similarity) function:

Do Perform feature Scaling before using the Gaussian Kernel

Landmarks are training examples as well

$$\|x - l\|^2 \rightarrow v = x - l$$

$$\begin{aligned} \|v\|^2 &= V_1^2 + V_2^2 + \dots + V_m^2 \\ &= (X_1 - L_1)^2 + (X_2 - L_2)^2 + \dots + (X_m - L_m)^2 \end{aligned}$$

1000 dimensions
1-5 dimensions

Other choices of Kernel

Note: Not all similarity functions $\text{similarity}(x, l)$ make valid Kernels. (Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM packages

Many off-the-shelf Kernels Available:

- Polynomial Kernel → the similarity is $K(x, l) = (x^T l)^d$, another one: $(x^T l)^3$, $(x^T l + 1)^3$, $(x^T l + 5)^4$

the general form is: $(x^T l + \text{constant})^d$

Almost Always usually performs worse than Gaussian Kernel

$x^T l$ Non Negatives

- More esoteric: String Kernel, Chi-Squared Kernel, histogram intersection Kernel, ...

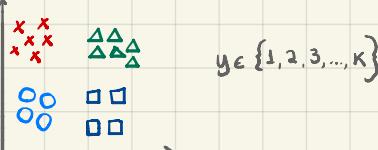
if your input data is text strings

optimizations run correctly, and do not diverge.
Buckets
Similarity between objects

If you are trying to decide among a few different choices of Kernel and are also choosing parameters such as C, σ^2 , etc.

You should choose whatever performs best on the cross-validation data

Multi-class classification



Many SVM packages already have built-in multi-class classification functionality. Otherwise, use one-vs-all method. (Train K SVMs, one to distinguish $y=i$ from the rest for $i: 1, 2, \dots, K$), get $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$ → pick class i with largest $(\theta^{(i)})^T x$.

Logistic regression vs SVM's

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

If m is large (relative to n): (E.g. $n=2000, m=10,000, m=10, \dots, 1000$)

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

If m is small, m is intermediate: ($m=1-2000, m=10-20,000$)

Use SVM with Gaussian Kernel

If m is small, m is large: ($m=1-2000, m=50,000+$)

Create/add more features, then use logistic regression or SVM without a Kernel

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

The optimization problem of SVM is convex → So, you don't need to worry about local optima. It always find a global minimum

or something close to.

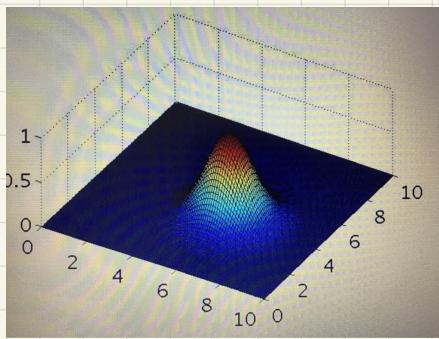


To find non-linear decision boundaries with the SVM, we need to first implement a Gaussian Kernel. You can think of the Gaussian Kernel as a similarity function that measures the "distance" between a pair of examples $(i^{(1)}, x^{(2)})$.

The Gaussian Kernel is also parameterized by a bandwidth parameter, σ , which determines how fast the similarity decreases (to 0) as the examples are further apart.

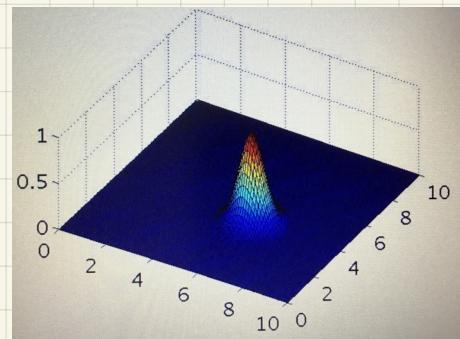
SVM Quiz → Important points

- If you have a decision boundary that is underfit to the training set, we'd like to lower bias/increase the variance of the SVM. We can do so by:
 - Either increasing the parameter C or decreasing σ^2 .
- When you are using a Gaussian Kernel, given by $\text{similarity}(x_i, l^{(m)}) = \exp\left(-\frac{\|x_i - l^{(m)}\|^2}{2\sigma^2}\right)$ and you have a plot of $f_2 = \text{similarity}(x_i, l^{(m)})$ when $\sigma^2 = 1$:



If you plot f_2
with $\sigma^2 = 0.25$

Then you have:



It shows a "narrower" Gaussian Kernel centered at the same location
Which is the effect of decreasing σ^2 .

Suppose you have a dataset with $m=10$ features and $m=5000$ examples.

After training your logistic regression classifier with gradient descent, you find that it has underfit the training set and does not achieve the desired performance on the training or cross validation sets. For this situation might be promising steps to take as follow:

- Use a SVM with Gaussian Kernel
 - by using a Gaussian Kernel, your model will have greater complexity and can avoid Underfitting the data.
- Create/add new polynomial features
 - When you add more features, you increase the variance of your model, reducing chances of underfitting

The maximum value of the Gaussian Kernel (i.e., $\text{sim}(x_i, l^{(m)})$) is 1.

When $x=l^{(m)}$, the Gaussian Kernel has value $\exp(0)=1$, and it is less than 1 otherwise.

It is important to perform feature normalization before using the Gaussian Kernel.

The similarity measure used by the Gaussian Kernel expects that the data lie in approximately the SAME range.

Week 8 - Unsupervised Learning

Clustering

• Introduction to unsupervised learning → No labels

• K-Means Algorithm: Is a method to automatically cluster similar data examples together. The intuition behind K-means is an iterative procedure that starts by guessing the initial centroids, and then refines this guess by repeatedly assigning examples to their closest centroids and then recomputing the centroids based on the assignments.

Input:

• K (number of clusters)

• Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ $x^{(i)} \in \mathbb{R}^n$ (drop $k_0=1$ convention)

Algorithm:

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

• Repeat {

cluster assignment step for $i=1$ to m
 $c^{(i)} := \text{index (from 1 to } K \text{) of cluster centroid closest to } x^{(i)}$ $\min_k \|x^{(i)} - \mu_k\|^2 \rightarrow c^{(i)}$

move centroids for $k=1$ to K

$\mu_k := \text{Average (mean) of points assigned to cluster } k$.

$$x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)} \rightarrow c^{(1)}=2, c^{(2)}=2, c^{(3)}=2, \dots, c^{(m)}=2$$

$$\mu_2 = \frac{1}{4} [x^{(1)} + x^{(2)} + x^{(3)} + \dots + x^{(m)}] \in \mathbb{R}^n$$

• Optimization Objective

• $c^{(i)}$ = index of cluster ($1, 2, \dots, K$) to which example $x^{(i)}$ is currently assigned

• μ_k = cluster centroid $\in (\mathbb{R}^n)$

$\mu_{c^{(i)}}$ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned $\rightarrow x^{(i)} \rightarrow c^{(i)} = 5 \rightarrow \mu_{c^{(i)}} = \mu_5$

Optimization Objective:

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

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

Distortion function

This cost function is also called by:

first minimizes J in respect to variables $c^{(i)}$

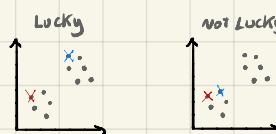
then minimizes J in respect to variables μ_k

• Random Initialization

→ Should have $K \leq m$

→ Randomly pick K training examples

→ Set μ_1, \dots, μ_K equal to these K examples



• Local Optima → You can initialize multiple times in order to assure that you'll get the best set.

For $i=1$ to 100 {

 Randomly initialize K-means

 Run K-means. Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_m$

 Compute Cost function/distortion $\rightarrow J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

}

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

(This is good) when you have a small number of clusters $K=2-10$

but even if have a large number of clusters it should give you a reasonable starting point to start from for finding a good set of clusters.

• Choosing the number of clusters

→ Elbow Method

If you run K-means using $K=3$ and $K=5$ and you find that the cost function J is much higher for $K=5$ than for $K=3$

↳ In the run with $K=5$, K-means got stuck in a bad local minimum. You should try re-running K-means with multiple random initializations.

→ Sometimes you are running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.

Dimensionality Reduction

↳ In this section we'll introduce PCA and show how it can be used for data compression to speed up learning algorithms as well as for visualizations of complex datasets

• Data Compression

· Reduce the data from 2D to 1D (e.g. cm/inches) or 3D to 2D, etc

· It makes some algorithms run much faster

• Data Visualization

• Principal Component Analysis (PCA) - Problem formulation

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

· Reduce from m -dimension to K dimension:

↳ find K vectors $u^{(1)}, u^{(2)}, \dots, u^{(K)}$ onto which to project the data, so as to minimize the projection error.

- PCA is NOT LINEAR regression

↳ PCA try to find a lower dimension surface onto which to project the data, so as to minimize this squared projection error.
In order to minimize the square distance between each point and the location of where it gets projected.

• PCA Algorithm

- Apply data preprocessing

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

Preprocessing / feature scaling / mean normalization:

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

Replace each $x_i^{(j)}$ with $x_i^{(j)} - \mu_j$ → then will have mean = 0

If different features on different scales (e.g., x_1 = size of house, x_2 = # of bedrooms), scale features to have comparable range of values.

Algorithm

Reduce data from n -dimensions to K -dimensions

Compute the "covariance matrix":

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

Compute "eigenvectors" of matrix Σ :

$$[U, S, V] = \text{SVD} / \text{Sigma} \rightarrow \text{Singular Value Decomposition (SVD)}$$

$$U = \begin{bmatrix} u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \end{bmatrix} \quad U \in \mathbb{R}^{n \times n} \\ u^{(1)}, \dots, u^{(n)}$$

$$U_{\text{reduce}} = \begin{bmatrix} u^{(1)} & u^{(2)} & \dots & u^{(K)} \end{bmatrix} \rightarrow Z = U_{\text{reduce}}^T \cdot X \rightarrow Z \in \mathbb{R}^{n \times K}$$

• SUMMARY

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

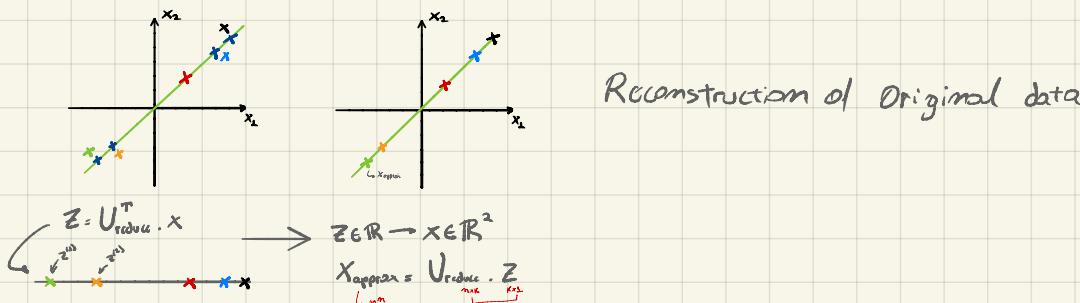
$$\text{Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T \quad \text{or the vectorized version Sigma} = (1/m) * X^T * X$$

$$[U, S, V] = \text{SVD} [\text{Sigma}]$$

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

$$Z = U_{\text{reduce}}^T * X$$

• Reconstruction from compressed Representation



• Choosing the Number of Principal Components

→ Average squared projection error: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$

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

Typically, choose K to be smallest value so that:

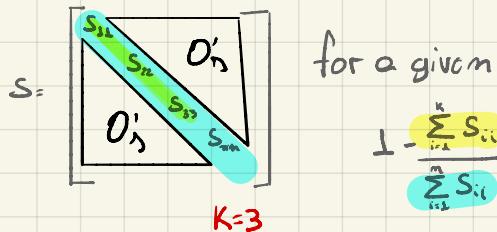
$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad (1\%) \rightarrow \text{It means "99\% of variance is retained"}$$

↳ Create a for loop for run PCA trying a range of K 's and computing $U_{\text{reduce}}, z^{(1)}, z^{(2)}, \dots, x_{\text{approx}}^{(1)}, \dots$, and checking if the % chosen of retained variance is satisfied

↳ Not efficient, because run PCA over and over

You can do it just iterating using S/diagonal matrix from SVD

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



for a given range of K pick the smallest value of K which:

$$1 - \frac{\sum_{i=1}^K S_{ii}}{\sum_{i=1}^m S_{ii}} \leq 0.01 \quad \text{OR} \quad \frac{\sum_{i=1}^K S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

Previously, we said that PCA chooses a direction $u^{(1)}$ for K directions $u^{(1)}, \dots, u^{(K)}$ onto which to project the data so as to minimize the squared projection error. Another way to say the same is that PCA tries to minimize: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - z^{(i)}\|^2$

• Advice for Applying PCA

Supervised learning speed up

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

Extract inputs: Unlabeled dataset: $x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{10000}$

↓ PCA

$$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^{1000}$$

$$\text{Now training set: } (z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA **ONLY** on the training set. This mapping can be applied as well to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test sets.

Applications of PCA

- Compression

- { Reduce memory/disk needed to store data
- Speed up learning algorithms
- Choose K by % of variance retain

- Visualization

- We can't plot over 3 dimensions, so PCA with $K=2$ or $K=3$ can help.

BAD USE OF PCA: To prevent overfitting

Use $z^{(i)}$ instead of $x^{(i)}$ to reduce the # of features to $K < n$, thus, fewer features, Less Likely to overfit **BAD**. Do not do this

This might work OK, but is not a good way to address overfitting.

Use **REGULARIZATION** instead

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

PCA is sometimes used where it shouldn't be

Design of ML system:

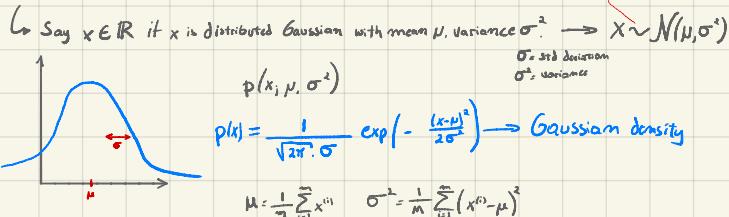
- Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$
- Train logistic regression $\{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\}$
- Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_{\theta}(z)$ on $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

→ How about doing the whole thing without using PCA?

Before implementing PCA, first try running whatever you want to do with the original/raw data $x^{(i)}$. Only if that doesn't do what you want, then implement PCA and consider using $z^{(i)}$.

Week 9 - Anomaly Detection

Gaussian Distribution or Normal Distribution



Anomaly detection: Algorithm

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

2- fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_m^2$

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

3: Given new example x , compute $p(x)$:

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

Anomaly if $p(x) < \epsilon$

Developing and evaluating an anomaly detection system

The importance of real-number evaluation

When developing a learning algorithm (choosing features, etc), making decisions is much easier if we have a way of evaluating our learning algorithm.

Assume we have some labeled data, of anomalous and non-anomalous examples ($y=0$ if normal, $y=1$ if anomalous).

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume just normal examples / not anomalous)

Cross Validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(n)}, y_{cv}^{(n)})$

Test set: $\underbrace{(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m)}, y_{test}^{(m)})}_{y=1}$

Aircraft engines example

\rightarrow 10,000 good (normal) engines

\rightarrow 20 flawed engines (anomalous)

Training set: 6000 good engines ($y=0$)

CV set: 2000 good engines ($y=0$), 10 anomalous ($y=1$)

Test set: 2000 good engines ($y=0$), 10 anomalous ($y=1$)

Algorithm Evaluation

Fit model $p(x)$ on Training set $\{x^{(1)}, \dots, x^{(m)}\}$

On a cross validation / test example x , predict:

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

Possible evaluation metrics:

\rightarrow Precision / Recall

$\rightarrow f_1$ Score

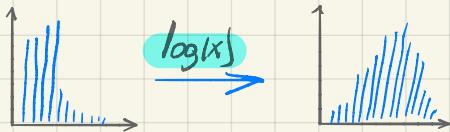
Can also use CV set to choose parameter ϵ .

Anomaly detection vs. Supervised learning

- Very small # of $y=1$
- Large # of $y=1, y=0$
- Large # of $y=0$
- Enough positive examples
- Many \neq types of anomalies
- for algorithm to get sense
- Hard for any algorithm
- of what positive examples are
- to learn from positive examples
- like, future positive examples
- likely to be similar to ones
- what the anomalies look like.
- in training set.
- Future anomalies may look different than the actual.

• Choosing what features to use

- ↳ It's better to have gaussian features
- ↳ In order to do it, we can try these:



Other options would be:

$$\begin{aligned} X_2 &\leftarrow \log(X_1 + 1) \\ X_3 &\leftarrow \sqrt{X_1^2} \\ X_4 &\leftarrow \sqrt[3]{X_1} \end{aligned}$$

- ↳ Choose features that might take on usually large or small values in the event of an anomaly.
- ↳ Sometimes you can try to divide one feature by another

• Multivariate Gaussian (Normal) distribution.

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

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

• Anomaly detection using the multivariate Gaussian distribution

Parameters: $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$

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

Parameter fitting:

Given training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

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

1- Fit the model $p(x)$ by setting:

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

2- Given a new example x , compute

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

Flag an anomaly if $p(x) < \epsilon$

Relationship to the Original model

The Multivariate Gaussian model is the same as original model when $\Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}$

Original model vs Multivariate Gaussian

$$p(x_1; \mu_1, \sigma_1^2) \times \dots \times p(x_m; \mu_m, \sigma_m^2)$$

- Manually create features to capture anomalies where x_1, x_2 take unusual combination of values. $X_3 = \frac{x_1}{x_2} = \text{CPU load} / \text{memory}$
- Computationally cheaper/scale better to large m , e.g. $m=10000$, $n=100000$
- OK even if m (training set size) is small

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

- Automatically captures correlations between features
- Computationally more expensive $(\frac{n^2}{2})$
- MUST have $m \geq n$ or else Σ is non-invertible
↳ $m \gg n \Rightarrow m \geq n$

Week 9 - Recommender System

They look at patterns of activities between different users and products to produce recommendations.

Problem formulation

$$n_u = \text{NO. users}$$

$$n_m = \text{NO. movies}$$

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

$$y^{(i,j)} = \text{rating given by user } j \text{ to movie } i \text{ (defined only if } r(i, j) = 1\text{)}$$

Content-based Recommendations

Problem formulation:

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

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

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

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

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

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

Optimization Objective:

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

$$\min_{\theta^{(j)}} \frac{1}{2} \sum_{i \in \text{rated}_j} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (\theta_k^{(j)})^2$$

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

$$\underbrace{\min_{\theta^{(1)}, \dots, \theta^{(n)}} \frac{1}{2} \sum_{j=1}^n \sum_{i \in \text{rated}_j} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n (\theta_k^{(j)})^2}_{J(\theta^{(1)}, \dots, \theta^{(n)})}$$

Gradient descent update:

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

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

Collaborative filtering

Given $x^{(1)}, \dots, x^{(m)}$ (and movie rating), can estimate $\theta^{(1)}, \dots, \theta^{(m)}$

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

Given $\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \dots$

Algorithm:

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

2. Minimize $J(x^{(1)}, \dots, x^{(m)}, \theta^{(1)}, \dots, \theta^{(m)})$ using gradient descent (or an advanced optimization algorithm). E.g.:

for every $j=1, \dots, n_u, i=1, \dots, n_m$:

$$x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j \in \text{rated}_i} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_j^{(i)} + \lambda x_k^{(i)} \right)$$

$$\theta_k^{(i)} := \theta_k^{(i)} - \alpha \left(\sum_{j \in \text{rated}_i} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x^{(i)} + \lambda \theta_k^{(i)} \right)$$

3. for a user with parameters θ and a movie with $(k \times m)$ features x , predict a star rating of $\theta^T x$.

Vectorization: Low rank matrix factorization

Movie	Alice (1)	Bob (2)	Carol (3)	Dane (4)
Laurel West	5	5	0	0
Romance Farmer	5	?	?	0
Cute puppies	?	4	0	?
abusive cat classes	0	0	5	4
Swords vs. Karate	0	0	5	?

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

$$\begin{aligned} &\text{Predicted ratings} \\ &[(\theta^{(1)})^T (x^{(1)}) \quad (\theta^{(1)})^T (x^{(2)}) \quad \dots \quad \dots \quad (\theta^{(1)})^T (x^{(n)})] \\ &[(\theta^{(2)})^T (x^{(1)}) \quad (\theta^{(2)})^T (x^{(2)}) \quad \dots \quad \dots \quad (\theta^{(2)})^T (x^{(n)})] \\ &\vdots \quad \vdots \\ &[(\theta^{(n)})^T (x^{(1)}) \quad \dots \quad \dots \quad \dots \quad (\theta^{(n)})^T (x^{(n)})] \end{aligned}$$

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

$$X \cdot \Theta^T$$

Low rank matrix factorization

• finding related movies

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

↳ $X_1 = \text{romance}$, $X_2 = \text{action}$, ...

→ How to find movies j related to movie i ?

↳ Small $\|x^{(i)} - x^{(j)}\|$ → movie j and i are "similar"

5 most similar movies to movie i :

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

• Mean Normalization

↳ for user j , on movie i , predict: $(\theta^{(i)})^T (X^{(j)}) + \beta_i$

→ this is the mean of the rates for movie i

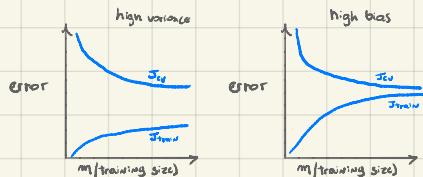
Week 10 - Large Scale Machine Learning

↳ Knowing how to handle big data is one of the most sought after skills in most companies.

• Learning with large datasets

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

↳ What if you have $m = 1,000,000$?



• Stochastic Gradient descent

Batch vs entire training set

Batch Gradient Descent

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

Repeat {

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

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

}

Stochastic Gradient Descent

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

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

1. Randomly shuffle dataset

2. Repeat {

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

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

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

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

$\rightarrow (x^{(i)}, y^{(i)}), (x^{(i)}, y^{(i)}), \dots$

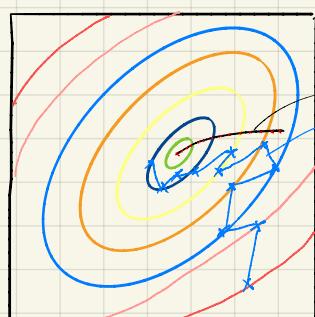
}

}

• Algorithm

1. Randomly shuffle (reorder) training

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



→ Batch Gradient descent

→ Stochastic Gradient descent → It "falls" to find a global minimum, but it finds its zone of global minimum, which isn't a good option.

• Mini-batch Gradient descent

→ Batch Gradient descent = Use all m examples in each iteration.

→ Stochastic Gradient descent = Use 1 example in each iteration.

→ Mini-batch Gradient descent = Use " b " examples in each iteration.

↳ Usually $b = 10$

- Algorithm: Mini-batch gradient descent

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

→ It's possible to apply vectorization.

Repeat { for $i=1, 11, 21, 31, \dots, 931$ {

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=1}^{10} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

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

Stochastic Gradient Descent - Convergence

- Checking for convergence:

→ Batch Gradient Descent:

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

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

→ Stochastic Gradient Descent

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

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

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

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

If the cost is increasing you can try to decrease the learning rate.

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

Online Learning

- Shipping service
- Product Search (Learning to search) → CTR
- Choosing special offers to show user
- Customized selection of news articles
- Product recommendation

Map-reduce and Data parallelism

Map-reduce → Batch Gradient Descent

$$\theta_j := \theta_j - \alpha \frac{1}{1000} \sum_{i=1}^{1000} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Training set

Instead of Machines it could be cores

- Machine 1 → Use $(x^{(1)}, y^{(1)}), \dots, (x^{(100)}, y^{(100)})$ → $\text{Temp}_j^{(1)} = \sum_{i=1}^{100} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$
- Machine 2 → Use $(x^{(1)}, y^{(1)}), \dots, (x^{(100)}, y^{(100)})$ → $\text{Temp}_j^{(2)} = \sum_{i=1}^{100} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$
- Machine 3 → Use $(x^{(1)}, y^{(1)}), \dots, (x^{(100)}, y^{(100)})$ → $\text{Temp}_j^{(3)} = \sum_{i=1}^{100} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$
- Machine 4 → Use $(x^{(1)}, y^{(1)}), \dots, (x^{(100)}, y^{(100)})$ → $\text{Temp}_j^{(4)} = \sum_{i=1}^{100} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$

$$\text{Combine: } \theta_j := \theta_j - \alpha \frac{1}{400} (\text{Temp}_j^{(1)} + \text{Temp}_j^{(2)} + \text{Temp}_j^{(3)} + \text{Temp}_j^{(4)})$$

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

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

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

$$\begin{aligned} J_{\text{train}}(\theta) &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) - (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \\ \frac{\partial}{\partial \theta_j} J_{\text{train}}(\theta) &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \end{aligned}$$

Week 11 - Application Example - Photo OCR

- Problem Description and pipeline:

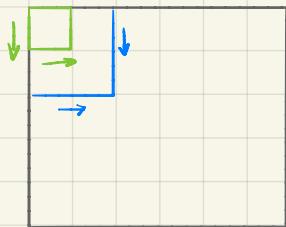
OCR = Optical Character Recognition



Pipeline



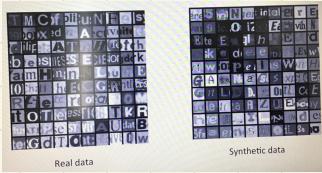
Sliding Window Detection



Step-size / Stride
It is used to "scan" the image with a fixed window that slides in the image.

Getting lots of Data: Artificial Data Synthesis

OCR:



this synthetic data can be obtained using font styles + Noise/distortions

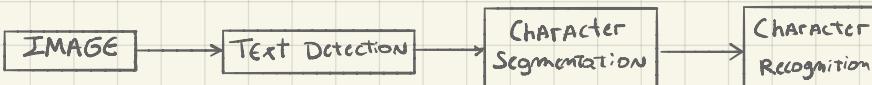
The synthetic data can be produced for a lot of applications that needs more training examples
↳ E.g. speech recognition (can use the original data with added background noises or bad cell phone connection)

The distortion introduced in your synthetic data should be representation of the type of noise/distortions in the test set.
Usually does not help to add purely random/meaningless noise to your data.

Discussion on Getting more data

1. Make sure you have a **LOW BIAS** classifier before expending the effort. (Plot curves). E.g. Keep increasing the number of features/number of hidden units in neural networks 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 → # hours?
 - "Crowd source" (E.g. Amazon Mechanical Turk)

Ceiling Analysis: What part of the pipeline to work on next

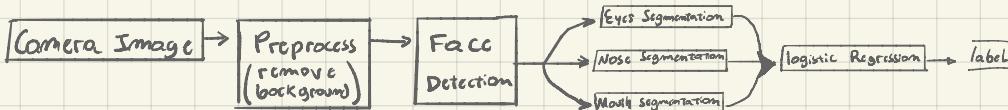


What part of the pipeline should you spend the most time trying to improve?

Truth table → Component	Accuracy
Overall System	72%
Text detection	89%
Character Segmentation	90%
Character Recognition	100%

+1% pp
+1% pp
+10% pp

face recognition from Images



If we got 100% of accuracy in each component, one per time, how would this affect the overall accuracy of the system?