

Supervised Learning

Linear regression

Step 1. Hypothesis:

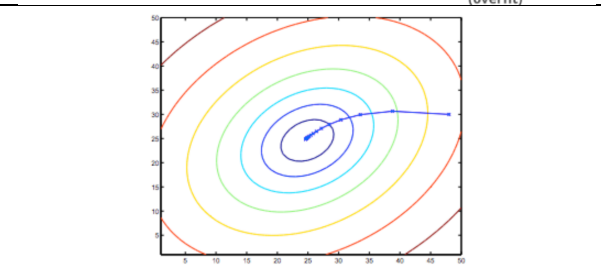
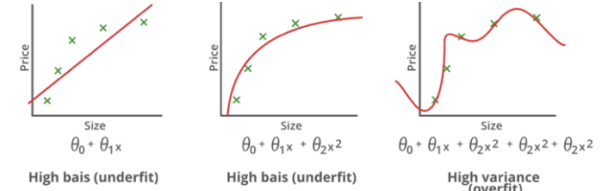
$$h_{\theta}(x)$$

Step 2. Cost

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

Step 3: Gradients

$$\begin{cases} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}, j = 1, 2, 3, \dots, n \end{cases}$$



Logistic regression

Step 1. Hypothesis:

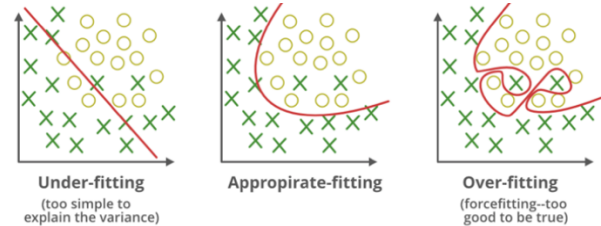
$$\begin{cases} h_{\theta}(z) = g(\theta^T x) \\ z = \theta^T x \\ g(z) = \frac{1}{1 + e^{-z}} \end{cases}$$

Step 2. Cost

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

Step 3. Gradients:

$$\begin{cases} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}, j = 1, 2, 3, \dots, n \end{cases}$$

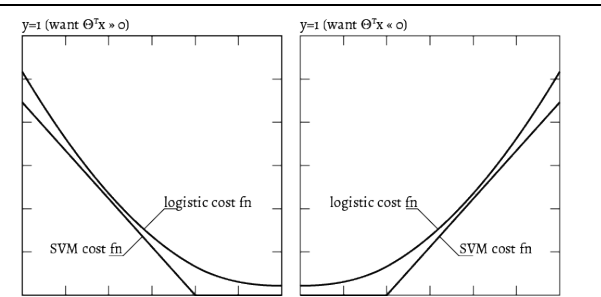


Support vector machines (SVM)

Cost

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

$$y^{(i)} = \begin{cases} 1, & \theta^T x^{(i)} \geq 1 \\ 0, & \theta^T x^{(i)} \leq -1 \end{cases}$$



SVM with Gaussian Kernel

Step 1. Hypothesis

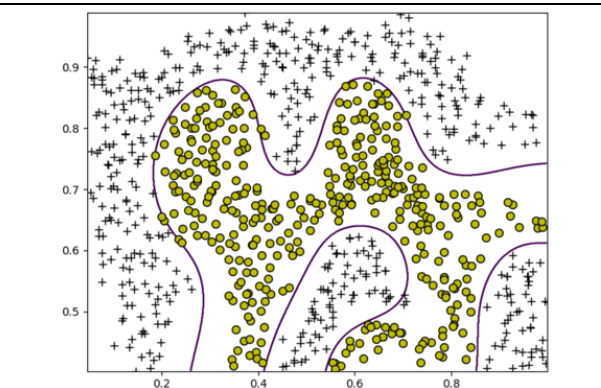
Given x , compute features $f \in \mathbb{R}^{m+1}$, parameters $\theta \in \mathbb{R}^{m+1}$
Predict "y=1" if $\theta^T f \geq 0$, $\theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m \geq 0$

Step 2. Training

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

$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right), \text{ or } = \left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

Predict "y = 1" if $\theta^T f_i \geq 0$



Neural network
Classification

Step 1. Randomly initialize weights

Initialize parameters $\Theta^{(1)}, \Theta^{(2)}, \dots, \Theta^{(L-1)}$ $[-\epsilon, \epsilon]$ (i.e. $-\epsilon \leq \theta_{ji}^{(l)} \leq \epsilon$)

Step 2. Forward propagation

$$h_{\Theta}(x^{(i)}) \in \mathbb{R}^K \quad (h_{\Theta}(X))_i = i^{th} \text{ output}$$

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

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

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

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

$$z^{(4)} = \Theta^{(3)} a^{(3)}$$

$$a^{(4)} = h_{\Theta}(x) = g(z^{(4)})$$

Step 3. Cost function $J(\Theta)$

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log(h_{\Theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_k) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{ji}^{(l)})^2$$

Step 4. Backpropagation to compute partial derivatives $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\Theta)$

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

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} * g'(z^{(2)})$$

$$g'(z^{(2)}) = a^{(2)} * (1 - a^{(2)})$$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} * g'(z^{(3)})$$

$$g'(z^{(3)}) = a^{(3)} * (1 - a^{(3)})$$

$$\delta^{(4)} = a^{(4)} - y$$

$$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

$$\left. \begin{aligned} D_{ij}^{(l)} &:= \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}, \text{ if } j \neq 0 \\ D_{ij}^{(l)} &:= \frac{1}{m} \Delta_{ij}^{(l)}, \text{ if } j = 0 \end{aligned} \right| D_{ij}^{(l)} = \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\Theta) = \frac{\partial J(\Theta)}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial \theta}$$

Step 5. Use gradient checking to compare $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\Theta)$ computed using

backpropagation vs. using numerical estimate of gradient of $J(\Theta)$.

Step 6. Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\Theta)$ as a function of parameters Θ .

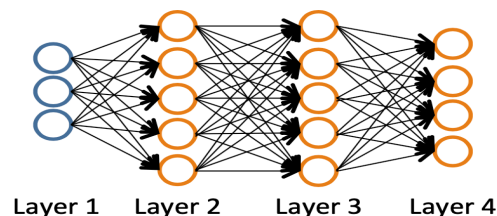
result = minimize(cost_func, initial_nn_params, method='CG',

jac=grad_func,

options={'disp': True, 'maxiter': 50.0})

nn_params = result.x

Jcost = result.fun



Unsupervised Learning

K-means

Step 1. Centroids

$c^{(i)} = \text{index of } \min \|x^{(i)} - \mu_j\|^2$
 $c^{(i)} \in \mathbb{R}^K, i = 1, 2, \dots, m$ denotes the index of cluster centroids closet to $x^{(i)}$

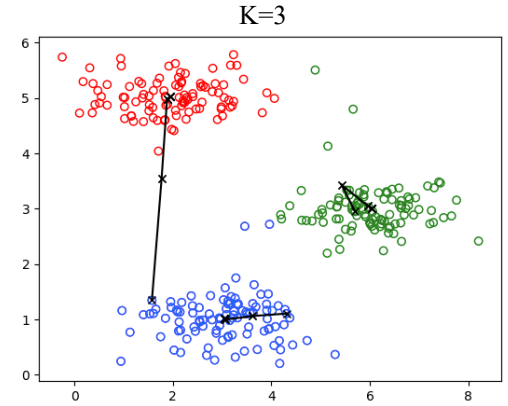
Step 2. Means

$$\mu_k = \frac{\sum_{i=1, \{c^{(i)}=k\}}^m x^{(i)}}{\sum_{i=1, \{c^{(i)}=k\}}^m 1}$$

$\mu_k \in \mathbb{R}^K, k = 1, 2, \dots, K$ denotes the average(mean) of points assigned to cluster k

Step 3. Cost function

$$J_{(c, \mu)} = \sum_i^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$



Principal Component Analysis (PCA)

Dimensionality Reduction

Data compression

Step 1. Feature scaling (Mean normalization)

Mean: $\bar{X} = \mu_j = \frac{1}{m} \sum_{i=1}^m x^{(i)}$

Standard deviation: $s = \sigma = \sqrt{\frac{1}{m-1} \sum_{i=1}^m (x^{(i)} - \mu)^2}$

Mean normalize: $x^{(i)} = \frac{x^{(i)} - \mu}{\sigma}$

Step 2. Calculate U, S, V.

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

U, S, V = numpy.linalg.svd(sigma)

$$U = \begin{pmatrix} | & | & & | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} & \dots & u^{(n)} \\ | & | & & | & | & | \end{pmatrix}$$

Ureduce = U[:, 0:K].T

Z = Ureduce * X = X_norm * U[:, 0:K]

X_approximate = X_recovered = Z * U[:, 0:K].T

Step 3. Pick the smallest value of k,

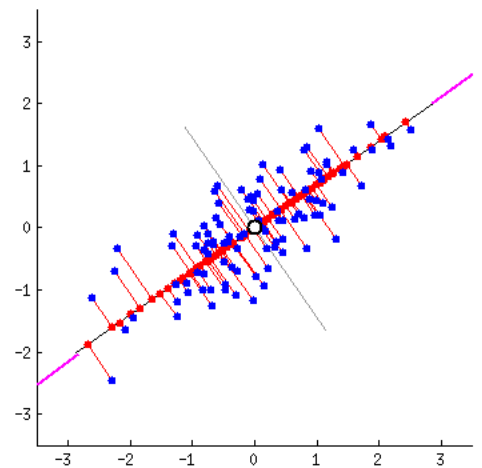
$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01?$$

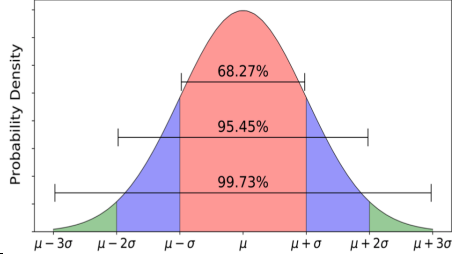
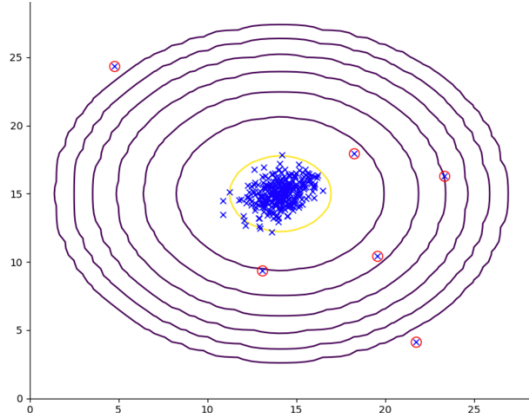
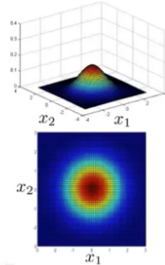
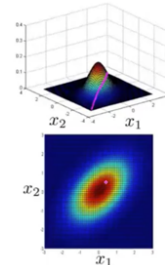
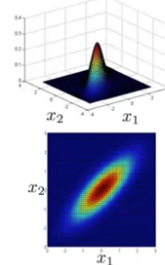
$$S = \begin{pmatrix} S_{11} & \dots & 0 \\ & S_{22} & \dots & 0 \\ \vdots & \vdots & S_{33} & \vdots & \vdots \\ & 0 & \dots & \ddots & \\ 0 & \dots & \dots & & S_{nn} \end{pmatrix}$$

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01 \rightarrow \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

99% of variance retained

2D → 1D



Anomaly Detection <i>Fraud detection, Manufacturing, Monitoring computers in a data center</i>	Gaussian (Normal) distribution $X \sim N(\mu, \sigma^2)$ Mean: $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$ Variance: $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$ Probability: $p(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	
Original model	<p><u>Step 1. Choose feature</u> Training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}, x^{(i)} \in \mathbb{R}^n$ Density estimation: $x_j \sim N(\mu_j, \sigma_j^2), j = 1, 2, \dots, n$ Choose features x_i that might be indicative of anomalous examples.</p> <p><u>Step 2. Fit parameters</u></p> $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$ $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$ <p><u>Step 3. Given new example $x \in \mathbb{R}^n$, compute $p(x)$</u> <u>Probability</u></p> $p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$ $y = \begin{cases} 1, & \text{if } p(x) < \epsilon(\text{anomaly}) \\ 0, & \text{if } p(x) \geq \epsilon(\text{normal}) \end{cases}$	
Multivariate Gaussian	<p><u>Step 1. Choose feature</u> Training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}, x^{(i)} \in \mathbb{R}^n$ Density estimation: $x_j \sim N(\mu_j, \sigma_j^2), j = 1, 2, \dots, n$</p> <p><u>Step 2. Fit parameters</u> Parameters: $\mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix) Mean: $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$ Variance: $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$</p> <p><u>Step 3. Given new example x, compute $p(x)$</u></p> <p>Diagonal Sigma: $\Sigma = \begin{pmatrix} \Sigma_1 & \dots & 0 \\ & \Sigma_2 & \dots & 0 \\ \vdots & \vdots & \Sigma_3 & \vdots & \vdots \\ 0 & 0 & \dots & \ddots & \Sigma_n \end{pmatrix}$</p> <p>Probability:</p> $p(x_j; \mu_j, \Sigma) = \frac{1}{\sqrt{2\pi} \Sigma ^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-u)}$ $y = \begin{cases} 1, & \text{if } p(x) < \epsilon(\text{anomaly}) \\ 0, & \text{if } p(x) \geq \epsilon(\text{normal}) \end{cases}$	<p>Multivariate Gaussian (Normal) examples</p> <div style="display: flex; justify-content: space-around;"> <div style="text-align: center;"> $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  </div> <div style="text-align: center;"> $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$  </div> <div style="text-align: center;"> $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$  </div> </div>

Machine learning is the science of getting computers to learn, without being explicitly programmed.

1 Supervised learning

Already know what our correct output should look like, having the idea that there is a relationship between the input and the output

1.1 Linear regression

Regression, continuous, individual

1) Data & Hypothesis

Input data: $X(m, n)$; dataset: $x^{(i)}, i = 1, 2, \dots, m$; features: $x_j, j = 1, 2, \dots, n$;

parameters: $\theta(n, 1)$; actual output: $y(m, 1)$.

$$X_{(m \times n)} = \begin{pmatrix} x_1^{(1)} & \dots & x_j^{(1)} & \dots & x_n^{(1)} \\ \vdots & & \vdots & & \vdots \\ x_1^{(i)} & \dots & x_j^{(i)} & \dots & x_n^{(i)} \\ \vdots & & \vdots & & \vdots \\ x_1^{(m)} & \dots & x_j^{(m)} & \dots & x_n^{(m)} \end{pmatrix}, \quad \theta^T_{(n \times 1)} = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_j \\ \vdots \\ \theta_n \end{pmatrix}, \quad Y^T_{(m \times 1)} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(i)} \\ \vdots \\ y^{(m)} \end{pmatrix}$$

2) Hypothesis:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^{0,1,2} x_2^{0,1,2} + \dots + \theta_n \prod_{j=1}^n x_j^{0,1,2,\dots,n}$$

3) Cost function:

We can measure the accuracy of our hypothesis function by using a **cost function**. This takes an average difference (actually a fancier version of an average) of all the results of the hypothesis with inputs from x's and the actual output y's.

Minimize cost:

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

4) Gradients

When specifically applied to the case of linear regression, a new form of the gradient descent equation can be derived. We can substitute our actual cost function and our actual hypothesis function and modify the equation to:

Rate: α , regulation: λ

$$G(\theta) = \frac{\partial J(\theta)}{\partial \theta}$$

Repeat until convergence: {

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

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}, \quad j = 1, 2, 3, \dots, n$$

}

$\alpha \frac{\lambda}{m} < 1$. Intuitively you can see it as reducing the value of θ_j by some amount on every update.

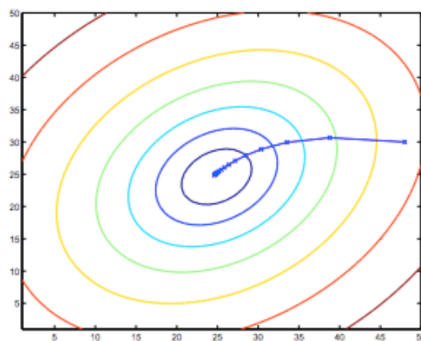


Figure 1 Example of batch gradient descent

Classification, Discrete, collective

1) Sigmoid function

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

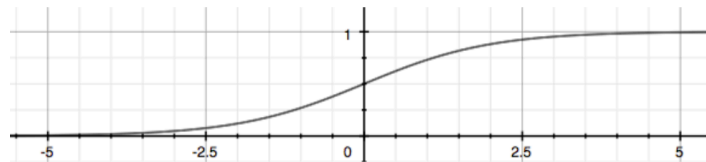
In order to get our discrete 0 or 1 classification, we can translate the output of the hypothesis function as follows:

$$\begin{cases} h_{\theta}(x) \geq 0.5 \Rightarrow \mathbf{y} = \mathbf{1} \\ h_{\theta}(x) < 0.5 \Rightarrow \mathbf{y} = \mathbf{0} \end{cases}$$

The way our logistic function g behaves is that when its input is greater than or equal to zero, its output is greater than or equal to 0.5:

$$g(z) \geq 0.5, \quad \text{when } z \geq 0$$

The following image shows us what the sigmoid function $g(z)$ looks like:



$$\begin{aligned} z = 0, e^0 = 1 &\Rightarrow g(z) = 1/2 \\ z \rightarrow \infty, e^{-\infty} \rightarrow 0 &\Rightarrow g(z) = 1 \\ z \rightarrow -\infty, e^{\infty} \rightarrow \infty &\Rightarrow g(z) = 0 \end{aligned}$$

So if our input to g is $\theta^T X$, then that means:

$$\begin{cases} h_{\theta}(z) = g(\theta^T x) \geq 0.5, \theta^T x \geq 0 \\ h_{\theta}(z) = g(\theta^T x) < 0.5, \theta^T x < 0 \end{cases}$$

From these statements we can now say:

$$\begin{cases} \theta^T x \geq 0 \Rightarrow \mathbf{y} = \mathbf{1} \\ \theta^T x < 0 \Rightarrow \mathbf{y} = \mathbf{0} \end{cases}$$

2) Cost function:

Minimize cost:

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

3) Gradients descent:

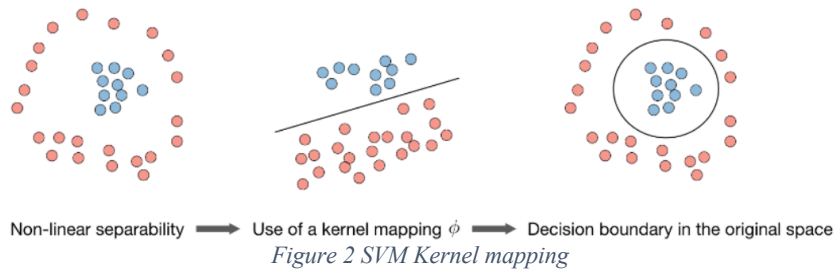
when computing the equation, we should continuously update the two following equations:

Repeat: {

$$\begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j &:= \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}, \quad j = 1, 2, 3, \dots, n \\ &\vdots \end{aligned}$$

1.3 Support vector machine (SVM)

SVM gives a cleaner, and more powerful way of learning complex non-linear functions.



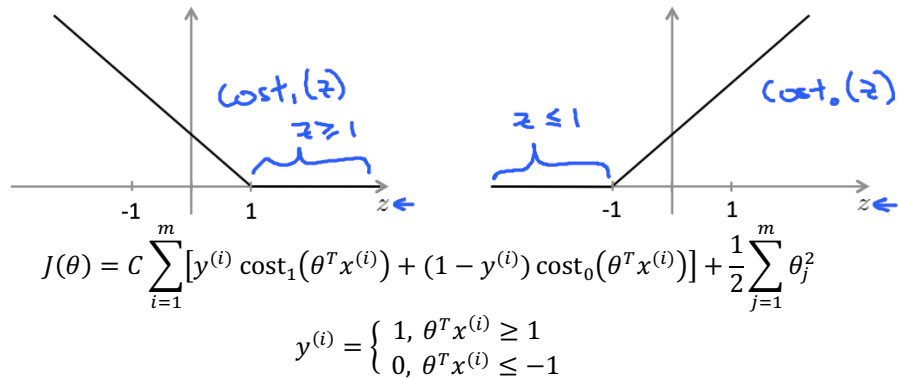
1.3.1 No kernel (“linear kernel”)

Predict “y = 1” if $\theta^T x > 0$

1) Hypothesis:

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

2) Cost function:



Large C: Lower bias, high variance.

Small C: Higher bias, low variance.

1.3.2 SVM with kernels, also called gaussian kernel

1) Hypothesis

Given x , compute features $f \in \mathbb{R}^{m+1}$, parameters $\theta \in \mathbb{R}^{m+1}$

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

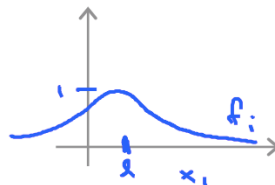
2) Training

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

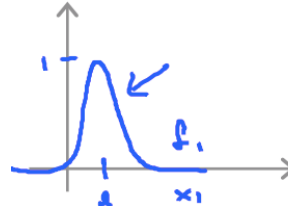
$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right), \text{ or } \left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

Predict “y = 1” if $\theta^T f_i \geq 0$

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.



3) Multiclass classification

Many SVM packages already have built-in multiclass 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

Logistic regression vs. SVMs

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

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

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

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

→ Use SVM with Gaussian kernel

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

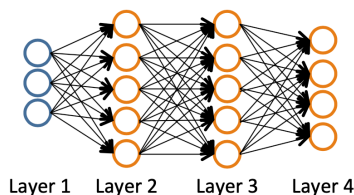
→ Create/add more features, then use logistic regression or SVM without a kernel

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



1.4 Neural Network

Non-linear Classification



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

L = total no. of layers in network

S_l = no. of units (not counting bias unit) in layer l

Binary classification

$y = 0$ or 1

1 output unit

$$h_{\theta} \in \mathbb{R}$$

$$S_L = 1 \quad (K = 1)$$

Multi-class classification (K classes)

$$y \in \mathbb{R}^K \quad \text{E.g.} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

pedestrian car motor truck

K output units

$$h_{\theta} \in \mathbb{R}^K$$

$$S_L = K \quad (K \geq 3)$$

- Training a neural network (e.g. $L=4$)
- Pick a network architecture (connectivity pattern between neurons)

- No. of input units: Dimension of features
- No. output units: Number of classes
- Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

1) Randomly initialize weights

Initialize parameters $\Theta^{(1)}, \Theta^{(2)}, \dots, \Theta^{(L-1)}$

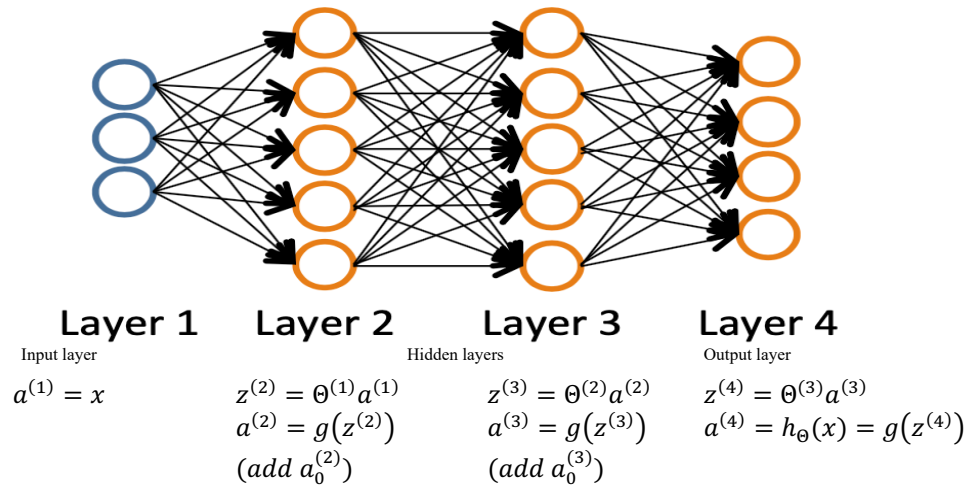
Initialize each $\Theta_{ji}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$ (i.e. $-\epsilon \leq \Theta_{ji}^{(l)} \leq \epsilon$)

epsilon = 0.12

theta = np.random.rand(input_layer, output_layer) * 2*epsilon - epsilon

2) Implement forward propagation to get $h_{\Theta}(x^{(i)})$ for any $x^{(i)}$

$$h_{\Theta}(x^{(i)}) \in \mathbb{R}^K \quad (h_{\Theta}(X))_i = i^{th} \text{ output}$$

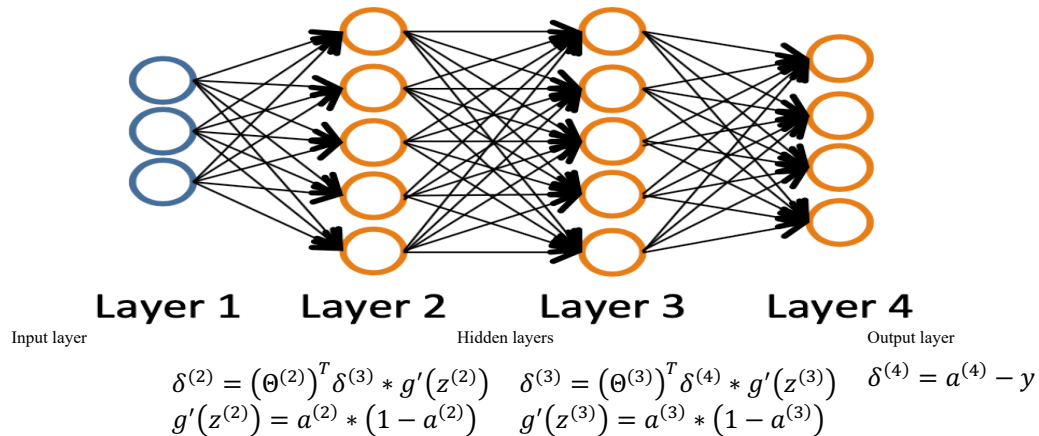


3) Implement code to compute cost function $J(\Theta)$

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log(h_{\Theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$

4) Implement backpropagation to compute partial derivatives $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$

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



Perform forward propagation and backpropagation using example $(x^{(i)}, y^{(i)})$
(Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l = 2, \dots, L$).

- Set $\Delta_{ij}^{(l)} = 0$ (from all, l, i, j).
- Set $a^{(1)} = x$
- Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \dots, L$
- Using y , compute $\delta^{(L)} = a^{(L)} - y$
- Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

$$\left. \begin{aligned} D_{ij}^{(l)} &:= \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}, \text{ if } j \neq 0 \\ D_{ij}^{(l)} &:= \frac{1}{m} \Delta_{ij}^{(l)}, \text{ if } j = 0 \end{aligned} \right| D_{ij}^{(l)} = \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = \frac{\partial J(\theta)}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial \theta}$$

5) Use gradient checking to compare $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$ computed using backpropagation vs. using numerical estimate of gradient of $J(\theta)$.

- i. Implement backpropagation to compute DELTA VECTOR (unrolled $D^{(1)}, D^{(2)}, \dots$).
- ii. Implement numerical gradient check to compute **gradient approximation**.
- iii. Make sure they give similar values.

```
diff = slin.norm(numgrad-grad)/slin.norm(numgrad+grad)
print('If your backpropagation implementation is correct, then \n\
the relative difference will be small (less than 1e-9). \n\
\nRelative Difference: ', diff)
```

iv. Turn off gradient checking. Using backprop code for learning.

- **Important:** Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of cost function your code will be very slow).

6) Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\theta)$ as a function of parameters θ .

```
result = minimize(cost_func, initial_nn_params, method='CG', jac=grad_func,
options={'disp': True, 'maxiter': 50.0})
nn_params = result.x
Jcost = result.fun
```

2 Unsupervised learning

Find hidden pattern in unlabeled data with little or no idea what our results should look like.
Don't necessarily know the effect of the variables.

2.1 K-means

1) Choose the value of K

Sometimes, you're 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.

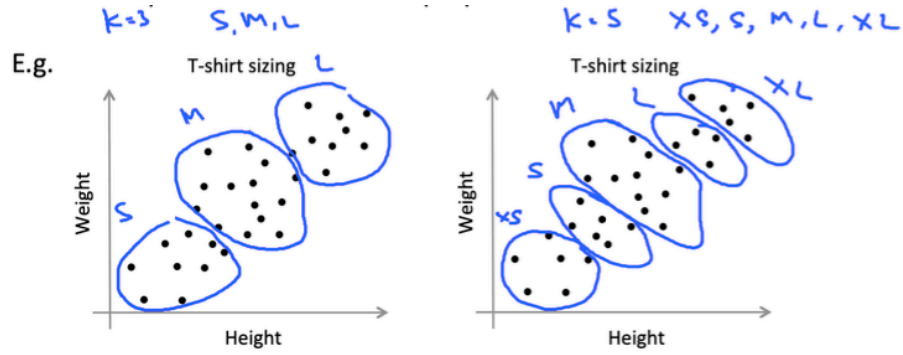


Figure 3 T-shirt size

2) Initialize centroids

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

For $i = 1$ to 100 {

Randomly pick examples from given points as K ($K < m$) different centroids

`initial_centroids = random.sample(X.tolist(), K)`

cost function

$$J_{(c, \mu)} = \sum_i^m \|x^{(i)} - \mu_{c(i)}\|^2$$

Return centroids of the smallest J .

}

3) Iteration

Set Iterate times, max_iters

For iters = 1 to max_iters {

$$c^{(i)} = \text{index of } \min \|x^{(i)} - \mu_j\|^2$$

$c^{(i)} \in \mathbb{R}^K, i = 1, 2, \dots, m$ denotes the index of cluster centroids closet to $x^{(i)}$.

$$\mu_k = \frac{\sum_{i=1, \{c^{(i)}=k\}}^m x^{(i)}}{\sum_{i=1, \{c^{(i)}=k\}}^m 1}$$

$\mu_k \in \mathbb{R}^K, k = 1, 2, \dots, K$ denotes the average(mean) of points assigned to cluster k .

}

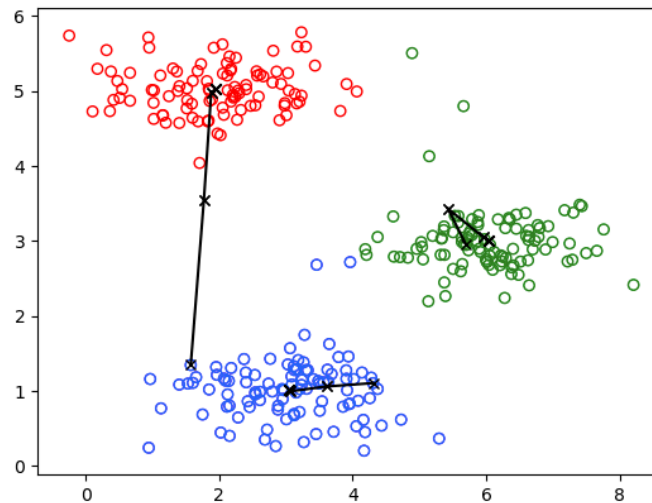


Figure 4 K clusters

2.2 Principal component analysis (PCA)

Dimensionality Reduction

Motivation:

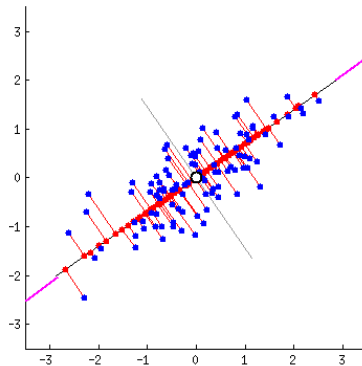


Figure 5 2 dimensions to 1 dimension

1) Data processing

Training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}, x^{(i)} \in \mathbb{R}^n$

Processing: feature scaling (mean normalization) to ensure every feature has zero mean.

i. Mean value:

$$\bar{X} = \mu_j = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

ii. Standard deviation

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

iii. Replace each $x^{(i)}$

$$x^{(i)} = \frac{x^{(i)} - \mu}{\sigma}$$

2) PCA algorithm.

Reduce data from n-dimensions to k-dimensions

$$\text{sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T = \frac{1}{m} X^T X \quad U = \begin{pmatrix} | & | & | & | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} & \dots & u^{(n)} \\ | & | & | & | & | & | \end{pmatrix}$$

Principal Component Analysis (PCA) algorithm:

```
U, S, V = numpy.linalg.svd(sigma)
```

Dimension matrix:

```
Ureduce = U[:, 0:K].T
```

Reduced K-dimensions data:

```
Z = Ureduce * X = X_norm * U[:, 0:K]
```

Recover data to n-dimensions:

```
X approximate = X recovered = Z * U[:, 0:K].T
```

3) **Choosing K** (number of principal components):

Pick the smallest value of K.

$$S = \begin{pmatrix} S_{11} & \dots & \dots & \mathbf{0} \\ \vdots & S_{22} & \dots & \mathbf{0} \\ \vdots & \vdots & S_{33} & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & S_{nn} \end{pmatrix}$$

Check if

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01? \quad 1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01 \Rightarrow \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

99% of variance retained.

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 on $\{(\cancel{z^{(1)}}), y^{(1)}), \dots, (\cancel{z^{(n)}}), y^{(m)})\}$
- - 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 $z^{(i)}$.

2.3 Anomaly detection

2.3.1 Original model

1) Choose feature:

Training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}, x^{(i)} \in \mathbb{R}^n$

Density estimation: $x_j \sim N(\mu_j, \sigma_j^2), j = 1, 2, \dots, n$

Choose features x_i that might be indicative of anomalous examples.

2) Fit parameters:

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

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

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

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

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

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

2.3.2 Multivariate Gaussian

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

1) Choose feature

Training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}, x^{(i)} \in \mathbb{R}^n$

Density estimation: $x_j \sim N(\mu_j, \sigma_j^2), j = 1, 2, \dots, n$

2) Fit parameters

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

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

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

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

$$p(x_j; \mu_j, \Sigma) = \frac{1}{\sqrt{2\pi} |\Sigma|^{\frac{1}{2}}} e^{\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)}$$

$$\Sigma = \begin{pmatrix} \Sigma_1 & \dots & \dots & \mathbf{0} \\ & \Sigma_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \Sigma_3 & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \ddots & \Sigma_n \end{pmatrix}$$

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

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

- The importance of real-number evaluation

I. 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 60%: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples/not anomalous)

Cross validation set 20%: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

Test set 20%: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

II. Algorithm evaluation

Possible evaluation metrics

Precision

$$\text{Precision}(P) = \frac{TP}{TP + FP}$$

Recall

$$\text{Recall}(R) = \frac{TP}{TP + FN}$$

F1-Score

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

Comparison between these two models

→ Original model

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

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

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

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

OK even if m (training set size) is small

vs. → Multivariate Gaussian

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

→ Automatically captures correlations between features

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

$$\Sigma^{-1}$$

Computationally more expensive

$$\rightarrow \Sigma \sim \frac{n^2}{2}$$

Must have $m > n$ or else Σ is non-invertible. $m \geq 10n$

Andrew Ng

2.3.3 Collaborative filtering

Recommender system

1) Problem formulation

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

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

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

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

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

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

2) Initialization

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

3) Cost function

- i. Minimize $\{x^{(1)}, \dots, x^{(n_m)}\}$ and $\{\theta^{(1)}, \dots, \theta^{(n_u)}\}$
Min $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$

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

- ii. using gradient descent (or an advanced optimization algorithm). E.g. for every $j = 1, 2, \dots, n_u$; $i = 1, 2, \dots, n_m$:
Gradients

$$\begin{cases} x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j:r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right) \theta_k^{(j)} + \lambda x_k^{(i)} \right) \\ \theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \end{cases}$$

4) Prediction.

For a user with parameters $\theta^{(j)}$ and a movie with (learned) features $x^{(i)}$, predict a star rating of $(\theta^{(j)})^T x^{(i)} + \mu_i$

μ_i is the mean of rated users for movie $x^{(i)}$.

$$\mu_i = \text{mean} \left(\sum_{j:r(i,j)=1}^n x_j^{(i)} \right)$$

Anomaly detection	vs.	Supervised learning
<ul style="list-style-type: none"> ➤ Very small number of positive examples (<u>$y = 1$</u>). (<u>0-20</u> is common). ➤ Large number of negative (<u>$y = 0$</u>) examples. <u>$p(x)$</u> ← ➤ <u>Many different “types” of anomalies</u>. Hard for any algorithm to learn from positive examples what the anomalies look like; ➤ future anomalies may look nothing like any of the anomalous examples we’ve seen so far. 		<ul style="list-style-type: none"> Large number of positive and negative examples. ← Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set. ← Spam ←

3 Tips

3.1 Model Selection and Train/Validation/Test Sets

One way to break down our dataset into the three sets is:

- Training set: 60%
- Cross validation set: 20%
- Test set: 20%

We can now calculate three separate error values for the three different sets using the following method:

- Optimize the parameters in Θ using the training set for each polynomial degree.
- Find the polynomial degree d with the least error using the cross-validation set.
- Estimate the generalization error using the test set with $J_{\text{test}}(\Theta(d))$, ($d = \text{theta from polynomial with lower error}$);

This way, the degree of the polynomial d has not been trained using the test set.

3.2 Gaussian (Normal) distribution

$$X \sim N(\mu, \sigma^2)$$

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

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

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

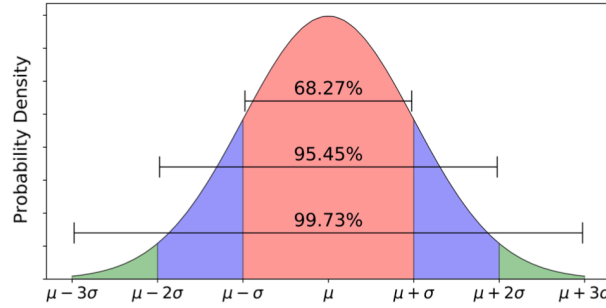


Figure 6 68% of the data is within 1 standard deviation, 95% is within 2 standard deviation, 99.7% is within 3 standard deviations

3.3 F₁-score

F-score is a measure of a test's accuracy.

		Actual class	
		1/+	0/-
Predicted class	1/+	True positive (TP)	False positive (FP)
	0/-	False negative (FN)	True negative (TN)

Precision

$$Precision(P) = \frac{TP}{\#Predicted P} = \frac{TP}{TP + FP}$$

Recall

$$Recall(R) = \frac{TP}{\#Actual P} = \frac{TP}{TP + FN}$$

F₁-score

$$F_1Score = 2 \frac{PR}{P + R}$$