

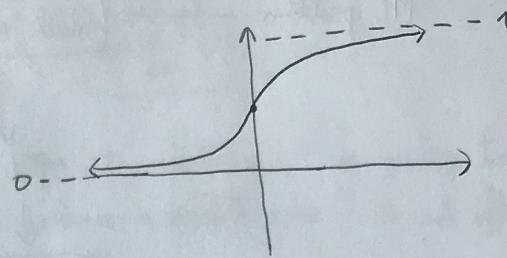
(24)

WEEK 7

FROM LR TO SVM:

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

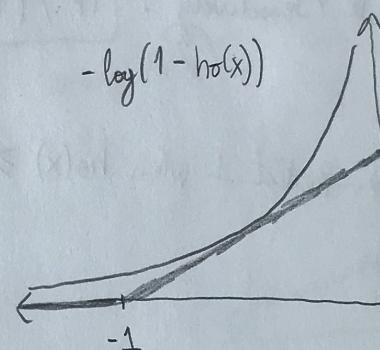
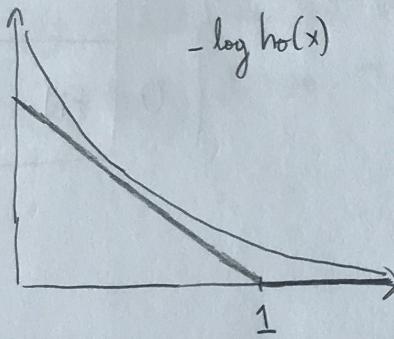
where $z = \theta^T x$



$$\begin{cases} y=1 \rightarrow h_{\theta}(x) \approx 1 \rightarrow \theta^T x \gg 0 \\ y=0 \rightarrow h_{\theta}(x) \approx 0 \rightarrow \theta^T x \ll 0 \end{cases}$$

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

$\underbrace{\text{cost}_1(\theta^T x^{(i)})}_{\text{For } y=1}$ $\underbrace{\text{cost}_0(\theta^T x^{(i)})}_{\text{For } y=0}$

For $y=1 \rightarrow$ we want $\theta^T x \gg 0$ For $y=0 \rightarrow$ we want $\theta^T x \ll 0$ 

$$\text{SVM} \rightarrow \min_{\theta} \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{cost}_1 + (1-y^{(i)}) \text{cost}_0 + \frac{\lambda}{2m} \sum_{j=0}^n \theta_j^2$$

A

B

$$\begin{cases} \text{LR} \rightarrow A + \lambda B \\ \text{SVM} \rightarrow CA + B \rightarrow C \text{ is like } \frac{1}{\lambda} \end{cases}$$

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

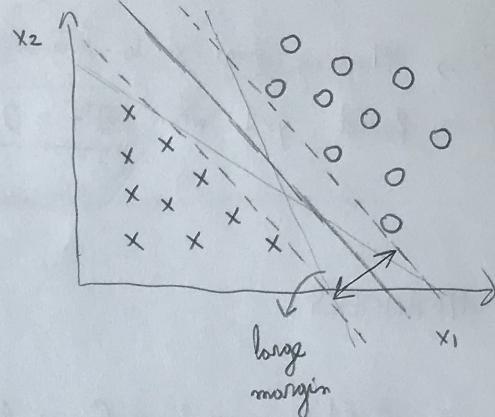
$$h_{\theta}(x) \begin{cases} 1 & \rightarrow \theta^T x \geq 0 \\ 0 & \rightarrow \text{otherwise} \end{cases}$$

(25)

LARGE MARGIN INTUITION :

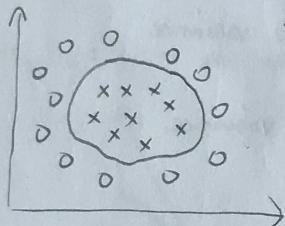
$$\begin{cases} y=1 \rightarrow \theta^T x \geq 1 \text{ (not just } \geq 0) \\ y=0 \rightarrow \theta^T x \leq -1 \text{ (not just } < 0) \end{cases}$$

↓
extra margin factor



→ C very large → sensitive to outliers

KERNELS :



$$\text{Predict } y=1 \rightarrow \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \theta_4 f_4 + \theta_5 f_5 + \dots \geq 0$$

↓ How to choose the features?

Ex: Proximity to landmarks $f^{(1)}, f^{(2)}, f^{(3)}$

$$\begin{cases} f_1 = \text{similarity}(x, l^{(1)}) = e^{-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}} \\ f_2 = \dots \\ f_3 = \dots \end{cases}$$

↳ kernels (gaussian kernel here)

When $x \approx l^{(1)} \rightarrow f_1 \approx 1$

When x far from $l^{(1)} \rightarrow f_1 \approx 0$

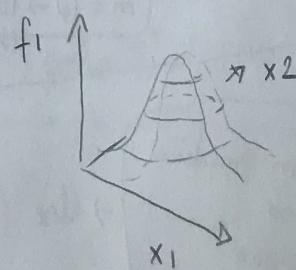
Plotting these similarity functions:

$\begin{cases} l^{(1)} \rightarrow \text{centre of the gaussian} \\ \sigma^2 \rightarrow \text{larger } \sigma^2 \text{ means smoother function} \end{cases}$

Ex: $\theta = \begin{bmatrix} -0.5 \\ 1 \\ 1 \\ 0 \end{bmatrix}$

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

→ x close to l^1 or $l^2 \rightarrow$ Predict 1
→ x far → Predict 0

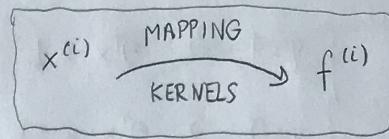


(26)

SVMs WITH KERNELS

→ Mapping from x to features f

→ Predict $y=1$ when $\theta^T f \geq 0$



PARAMETERS:

→ $C \left(\frac{1}{\lambda}\right)$

- ↳ Large $C \rightarrow$ lower bias / high variance
- ↳ Small $C \rightarrow$ higher bias / low variance

→ σ^2

- ↳ Large $\sigma^2 \rightarrow$ more smooth = higher bias / lower variance
- ↳ Small $\sigma^2 \rightarrow$ more abrupt = low bias / high variance

To avoid overfitting → pick small C and large σ^2

SVM OVERVIEW:

→ Specify C and kernel (if needed)

- ↳ No kernel (i.e. linear) $\rightarrow \theta^T x \geq 0$
- ↳ Kernel $\rightarrow \theta^T f \geq 0$

→ Use feature scaling before applying a kernel.

→ Multiclass SVM \rightarrow Train K SVMs if K classes, and pick largest one

$n \rightarrow$ features
 $m \rightarrow$ training samples

$n \geq m$
 $n = 10000$
 $m = 10 \rightarrow 1000$

→ Use LR or SVM without kernel

n small, m medium
 $n = 1 \rightarrow 1000$
 $m = 10 \rightarrow 10000$

→ Use SVM with gaussian kernel

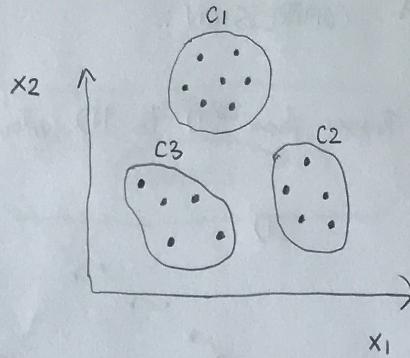
n small m large
 $n = 1 \rightarrow 1000$
 $m = 50000 +$

→ Add more features and use LR or SVM without kernel

→ K-MEANS CLUSTERING:

Input: → K (number of clusters)

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



① Random init. of K cluster centroids $\mu_1, \mu_2, \dots, \mu_K$

② Do [for each sample: assign $c^{(i)}$ → index of cluster centroid closest to $x^{(i)}$]

cluster assignment
move
centroids [for each cluster: compute $\mu_k \rightarrow$ average of points assigned to cluster k ,

OPT. OBJECTIVE
(DISTORTION)

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

→ RANDOM INIT:

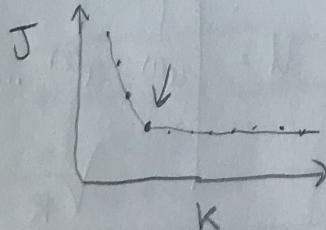
RUN ALGO MULTIPLE TIMES:

- ↳ Random. init. of μ s to training data pts
- ↳ Run algo and get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$
- ↳ Compute $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

Pick clustering that gave the lowest cost $J()$.

→ CHOOSING NUMBER OF CLUSTERS:

↳ Elbow method:

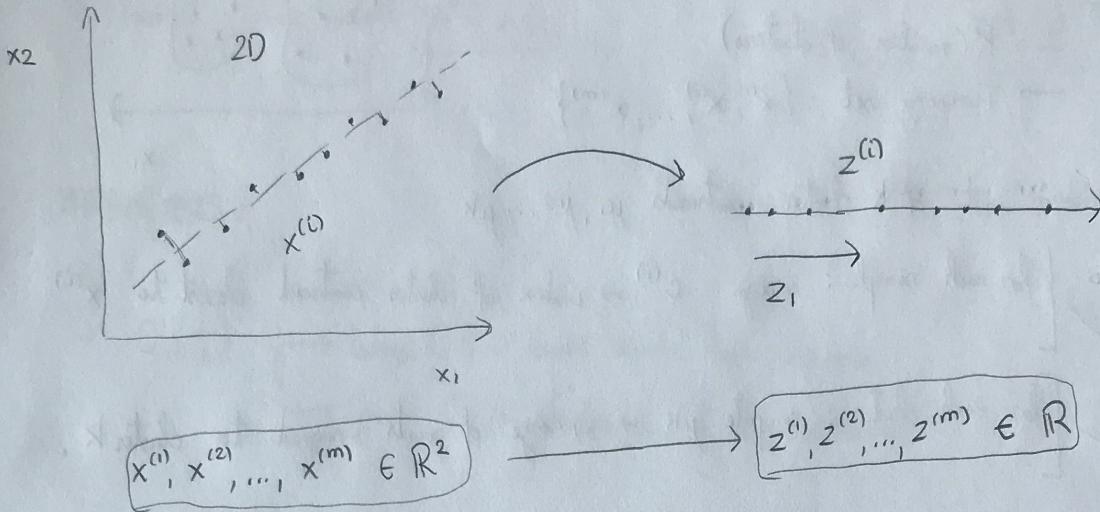


↳ Based on later purpose; we metric on how well it performs on a later task

(28)

DATA COMPRESSION :

→ Reduce from 2D to 1D when correlated features



→ 3D from 2D : From $\begin{bmatrix} x_1, x_2, x_3 \end{bmatrix} \in \mathbb{R}^3$ to $\begin{bmatrix} z_1, z_2 \end{bmatrix} \in \mathbb{R}^2$

DATA VISUALIZATION :

Ex: Reduce from 50D (50 features) to 2D (2 features).

PRINCIPAL COMPONENT ANALYSIS (PCA):

→ Reduce from n-dim to k-dim : find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data to minimize the projection error.

① Data processing:

TRAINING SET : $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

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

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

② ALGO: (from n-d to k-d)

↓ COVARIANCE MATRIX : $\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)})(x^{(i)})^T$ $(n \times n)$

COMPUTE EIGENVECTORS : $(n \times n)$

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

COMPUTE Z

$Z = U_k^T X$

(29)

Example 3D to 2D:

$$\{x_1, x_2, x_3\} \in \mathbb{R}^3$$

$$U_3 = \begin{bmatrix} U_{11} & U_{12} & U_{13} \\ U_{21} & U_{22} & U_{23} \\ \underbrace{U_{31} & U_{32} & U_{33}}_{3 \times 3} \end{bmatrix} \rightarrow U_2 = [U^{(1)} \ U^{(2)}] = \begin{bmatrix} U_{11} & U_{12} \\ U_{21} & U_{22} \\ U_{31} & U_{32} \end{bmatrix}_{3 \times 2}$$

$$Z = U_2^T x = \begin{bmatrix} U_{11} & U_{21} & U_{31} \\ U_{12} & U_{22} & U_{32} \end{bmatrix}_{2 \times 3} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{3 \times 1}}_{\text{3x1}} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}_{2 \times 1} \quad (2 \times 1)$$

RECONSTRUCTION :

$$Z = U_{\text{reduced}}^T \cdot x \rightarrow x_{\text{approx.}} = U_{\text{reduced}} \cdot Z$$

CHOOSING NUMBER OF PCs :

↳ Choose smallest value k so that

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

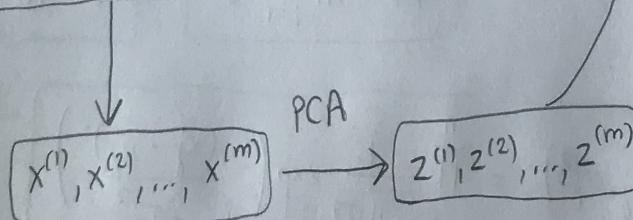
↳ total variation in the data

PCA ADVICE :

99% of variance is retained

SUPERVISED LEARNING
DATASET :
 $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

TRAIN ALGO ON:
 $(z^{(1)}, x^{(1)}), (z^{(2)}, x^{(2)}), \dots, (z^{(m)}, x^{(m)})$

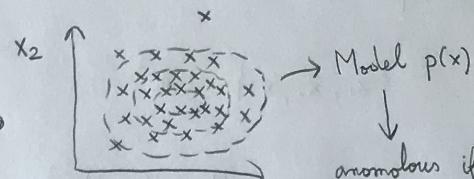


TEST ON:
 $x_{\text{CV}}^{(i)} \rightarrow z_{\text{CV}}^{(i)}$
 $x_{\text{test}}^{(i)} \rightarrow z_{\text{test}}^{(i)}$

(30)

WEEK (9)

ANOMALY DETECTION MOTIVATION:

2 features $\rightarrow x_1, x_2$ Dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \rightarrow$ density estimation x_{test} anomalous if $p(x_{\text{test}}) < \epsilon$ Is x_{test} anomalous?

GAUSSIAN (NORMAL) DISTRIBUTION:

$$x \sim \mathcal{N}(\mu, \sigma^2) \rightarrow p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

PARAM. ESTIMATION

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

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

DENSITY ESTIMATION:

$$\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \rightarrow \begin{cases} x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2) \\ x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2) \\ \vdots \\ x_n \sim \mathcal{N}(\mu_n, \sigma_n^2) \end{cases}$$

$$\rightarrow p(x) = p(x_1; \mu_1, \sigma_1^2) \cdot p(x_2; \mu_2, \sigma_2^2) \cdots p(x_n; \mu_n, \sigma_n^2)$$

ANOMALY DETECTION ALGO:

① Relevant features selection

② Fit params. $\mu_i \rightarrow \mu_n$ and $\sigma_i^2 \rightarrow \sigma_n^2$

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

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

③ Given new example x :

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

ANOMALY IF

$$p(x) < \epsilon$$

(31)

ANOMALY DETECTION DEV. AND. EVAL.

$$\begin{cases} y=0 \rightarrow \text{normal} \\ y=1 \rightarrow \text{anomalous} \end{cases}$$

TRAINING SET $\left\{ x^{(1)}, x^{(2)}, \dots, x^{(m)} \right\}$
Only normal examples

E.g. Training: 6000 normal
CV \rightarrow 2000 normal, 10 anom.
Test \rightarrow 2000 normal, 10 anom.

CV SET $\left\{ (x_{cv}^{(1)}, y_{cv}^{(1)}) \dots (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})}) \right\}$
 \hookrightarrow Some labelled $y=1$

TEST SET $\left\{ (x_t^{(1)}, y_t^{(1)}) \dots (x_t^{(m_t)}, y_t^{(m_t)}) \right\}$
 \hookrightarrow Some labelled $y=1$

① Fit $p(x)$ on the training set \rightarrow only normal samples

② On CV / test sample $x \rightarrow$ predict

$$\begin{cases} y=1 & \text{when } p(x) < \varepsilon \\ y=0 & \text{when } p(x) \geq \varepsilon \end{cases}$$

Use CV to choose value of ε

③ Use relevant metrics to evaluate algo:

$$\begin{cases} \rightarrow \text{TP} / \text{FP} / \text{FN} / \text{TN} \\ \rightarrow \text{Precision, recall} \\ \rightarrow F_1 \text{ score} \end{cases}$$

ANOMALY DETECTION VS SUPERVISED LEARNING :

\rightarrow Very small number of anomalous samples

\hookrightarrow estimate $p(x)$ from $y=0$

\rightarrow Many different types of anomalies

\rightarrow Future anomalies may be very different

\rightarrow Large number of both labels

\rightarrow Future positive samples likely to be similar to those in the training set.

FEATURE SELECTION:

→ If non gaussian: $\begin{cases} x_1 \rightarrow \log(x_1) \\ x_1 \rightarrow x_1^{1/3} \dots \end{cases}$

→ Choose new features that might be small/large for anomalies: $\begin{cases} x_3 \rightarrow \frac{x_1}{x_2} \\ x_3 \rightarrow \frac{x_1^2}{x_2} \end{cases}$

MULTIVARIATE GAUSSIAN DISTRIBUTION:

↳ When there may be a correlation between features

↳ Model $p(x)$ in one go with $\mu \in \mathbb{R}^n$ and $\Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

2 features

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

→ $\mu \rightarrow$ translation of centre of the distribution

→ diagonals in Σ : larger = smoother

→ other elements in Σ : correlation between x_1 and $x_2 \rightarrow$ elliptic

↓
ALGO.

① Fit model $p(x)$ wrong

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

② Given new sample x , compute:

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

③ Flag anomaly if $p(x) < \epsilon$

TYPES OF GRADIENT DESCENT ALGOS

① BATCH GD

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

REPEAT

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

for every $j = 0 \dots n$

② STOCHASTIC GD

Randomly shuffle training samples

REPEAT 1-10xFOR $i = 1, \dots, m$

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

FOR $j = 0 \dots n$ One iteration = all m samples

One iteration = only one sample

③ MINI-BATCH GD

mini batch size = b

$$1 \leq b \leq m$$

E.g. $\begin{cases} b = 10 \\ m = 100 \end{cases}$

LOOP $\left\{ \begin{array}{l} \text{for } i = 1, 11, 21, 31, \dots, 91 \\ \theta_j := \theta_j - \frac{\alpha}{10} \sum_{k=1}^{10} (h_{\theta}(x^{(k)}) - y^{(k)}) x_j^{(k)} \\ (\text{for every } j = 0, \dots, n) \end{array} \right.$

MAP - REDUCE

↓ Batch GD

- { Machine 1 → use samples 1 to 100 → $\text{temp}_j^{(1)} = \sum_{i=1}^{100} (\text{h}_0(x^{(i)}) - y^{(i)}) x_j^{(i)}$
- Machine 2 → 101 to 200 → $\text{temp}_j^{(2)}$
- Machine 3 → 201 to 300 → $\text{temp}_j^{(3)}$
- Machine 4 → 301 to 400 → $\text{temp}_j^{(4)}$

→ COMBINE ALL

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

for $j = 0 \rightarrow n$