# Alex Ivensky - ECE 1395 Midterm Sheet

God help me

## General Concepts:

Supervised Learning: given training data
$\rightarrow$ regression (continuous), classification (discrete)

Unsupervised Learning: finds hidden patterns
$\rightarrow$ clustering

## Linear Regression

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

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

(prediction, right answer)

Gradient Descent:

Repeat {
$\quad \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$
}  simultaneously update for every $j$

- Feature Scaling:

$x_{j,new} = \dfrac{x_j - \mu_j}{\sigma_j}$  ($\mu_j$ = mean $j$, $\sigma_j$ = std. dev. $j$)

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

## Logistic Regression (classification)

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

- $h_\theta(x) = g(\theta^T x)$.

$\rightarrow h(x) = P(y=1 | x)$

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

Want $\min_\theta J(\theta)$:

Repeat {
$\quad \theta_j := \theta_j - \alpha \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \cdot \frac{1}{m}$
}
(simultaneously update all $\theta_j$)

## Regularization
now also minimizing magnitudes of $\theta$ to reduce overfitting.

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

**Gradient descent**

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

$\quad \theta_j := \theta_j - \alpha \quad \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$  If $\lambda > 0$,
$\quad\quad (j = \cancel{0}, 1, 2, 3, \dots, n)$
}
$\theta_j := \theta_j(1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^{m}(h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$

$\theta = \left( X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix} \right)^{-1} X^T y$

## KNN

- Classification:

$y = \text{argmax}_c \#(y_i = c)$

$y = \text{argmax}_c \sum_{i \in N_K(x)} I(y_i = c)$  (set of K-NN's)

$\Pr(Y = j | X = x_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} I(y_i = j)$

- Regression:

$y = \frac{1}{K} \sum_{i \in N_K(x)} y_i$

## Weighted KNN

- Neighbors weighted differently:
  - Use all samples, i.e. K = N
  - Weight on i-th sample: $w_i = e^{\frac{-||x - x_i||^2}{\sigma^2}}$
  - σ = the bandwidth parameter, expresses how quickly our weight function "drops off" as points get further and further from the query x

- Classification:
$y = \text{argmax}_c \sum_{i=1}^{N} w_i I(y_i = c)$

- Regression:
$y = \frac{\sum_{i=1}^{N} w_i y_i}{\sum_{i=1}^{N} w_i}$

$d(\mathbf{x}, \mathbf{z}) = \left[ \sum_{i=1}^{D} (x_i - z_i)^2 \right]^{\frac{1}{2}}$

## Perceptron   linear classifier

Assume $y \in \{-1, 1\}$

$h_\theta(x) = f(\theta_T x)$ where $f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$

So $\quad \theta^T x > 0 \rightarrow$ classify **1**
$\quad\quad \theta^T x < 0 \rightarrow$ classify **-1**.

$J(\theta) = -\sum_{i \in M} \theta^T x^{(i)} y^{(i)}$    M is set of wrongly classified vectors.

## Perceptron Gradient Descent:

$\theta^{(\tau+1)} = \theta^{(\tau)} + \alpha \sum_{i \in M} x^{(i)} y^{(i)}$

## PCA   method of dimensionality reduction.

$C = (X - \mu)(X - \mu)^T$, covariance matrix

$= Q \Lambda Q^{-1}$

$= \begin{bmatrix} | & & | \\ ev_1 & \dots & ev_n \\ | & & | \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} - & ev_1^T & - \\ & \vdots & \\ - & ev_n^T & - \end{bmatrix}$

(eigenvectors) (eigenvalues)

% variance captured by the top k evs

$= \dfrac{\sum_{i=1}^{k} \lambda_i}{\sum_{j=1}^{n} \lambda_j}$    want $\geq 0.9$.

- Given: $M$ data points $x_1, .., x_M$ in $R^d$ where $d$ is big ($d \gg M$)

- We want some directions, $u_1, u_2, ..., u_k, ..., u_d$, in $R^d$ that capture most of the variation of the $x_i$. The coefficients would be: $\alpha_k = u_k^T (x_i - \mu)$
($\mu$: mean of data points)

Let $\Phi_i$ be the (very big vector of length $d$) that is face image $I$ with the mean image subtracted.
Define $C = \frac{1}{M}\sum \Phi_i \Phi_i^T = AA^T$
where $A = [\Phi_1 \Phi_2 ... \Phi_M]$ is the matrix of faces, and is $d \times M$.

*Note:* $C$ is a huge $d \times d$ matrix (remember $d$ is the length of the vector of the image).

## Recognition with eigenfaces

- Given novel image **x**:
  - Project onto subspace:
    $$w_1, ..., w_k = [u^T_1 (\mathbf{x} - \mu), ..., u^T_k (\mathbf{x} - \mu)]$$
  - Optional: check reconstruction error $x - \hat{x}$ to determine whether image is really a face
  - Classify as closest training face in k-dimensional subspace

## SVM.

$x_i$ positive ($y_i = 1$): $\quad \theta^T x_i \geq \boxed{1}$

$x_i$ negative ($y_i = -1$): $\quad \theta^T x_i \leq \boxed{-1}$

For support vector, $\quad \theta^T x_i = \pm 1$

- Want line that maximizes the margin.

$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$
$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

For support, vectors, $\quad \mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

Distance between point and line: $\quad \dfrac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$

**Margin**

For support vectors:
$$\frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|} = \frac{\pm 1}{\|\mathbf{w}\|} \qquad M = \left| \frac{1}{\|\mathbf{w}\|} - \frac{-1}{\|\mathbf{w}\|} \right| = \frac{2}{\|\mathbf{w}\|}$$

Minimize $\dfrac{1}{2}\mathbf{w}^T \mathbf{w}$

Subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

- Solution: $\theta = \sum_i \alpha_i y_i x_i$
$\forall$ s.v

Learned weight — Support vector

## The "kernel trick"

- The linear classifier relies on dot product between vectors $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$
- If every data point is mapped into high-dimensional space via some transformation $\Phi$: $\mathbf{x}_i \rightarrow \varphi(\mathbf{x}_i)$, the dot product becomes: $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$
- A *kernel function* is similarity function that corresponds to an inner product in some expanded feature space
- *The kernel trick*: instead of explicitly computing the lifting transformation $\varphi(\mathbf{x})$, define a kernel function K such that: $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$

## Using SVMs

$x_1, x_i, x_2,$  3rd feature → 150k  $x_2$ $x_3$ $x_1^2, x_2^2, y_3^2,$ $x$. 5000

1. Select a kernel function.
   This implicitly assumes a mapping to a higher dimensional (yet, not known) space.

2. Compute pairwise kernel values between labeled examples. g opt... £

3. Use this "kernel matrix" to solve for SVM support vectors & alpha weights.

4. To classify a new example: compute kernel values between new input and support vectors, apply alpha weights, check sign of output. $x_{t}$ $f(x_{t+}) = \sum N \cdot y \cdot K$

## Soft-margin SVMs (allow misclassification)

Misclassification cost  $N$ # data samples  Slack variable

$$\min_w \quad \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{N} \xi_i$$

The $w$ that minimizes...  Maximize margin  Minimize misclassification

$$\text{subject to} \quad y_i w^T x_i \geq 1 - \xi_i,$$
$$\xi_i \geq 0, \quad \forall i = 1, \ldots, N.$$

- Classification function:
$$f(x) = \text{sign}(\theta \cdot x)$$
$$= \text{sign}\left(\sum_i \alpha_i y_i x_i \cdot x\right)$$

If f(x) < 0, classify as negative, otherwise classify as positive.

- Notice that it relies on an *inner product* between the test point *x* and the support vectors $x_i$
- (Solving the optimization problem also involves computing the inner products $x_i \cdot x_j$ between all pairs of training points)

$$D = \frac{\theta^T x}{\|\theta\|}$$
L2-norm