1. Describe the regression task

- We have a domain $X \in \mathbb{R}^d$ and a label set $y \in \mathbb{R}$
- The hypothesis set is $\mathcal{H}_{reg} : \mathbb{R}^d \to \mathbb{R}$

OBJECTIVE: $\qquad h^* = \arg\min_{h \in \mathcal{H}_{reg}} L_d(h)$

2. Introduce the linear regression model class and derive the optimal solution

- Linear hypothesis class: $\mathcal{H}_{lin} = \{ x \to \langle w, x \rangle + b : w \in \mathbb{R}^d, b \in \mathbb{R} \}$

- Squared loss: $\ell_2(h, (x,y)) \overset{def}{=} (h(x) - y)^2$

ERM: $\qquad L_s(h) = \frac{1}{m} \sum_{i=1}^{m} (h(x_i) - y_i)^2$

Hom. COORD. $w^* = \arg\min_w \frac{1}{m} \sum_{i=1}^{m} (\langle w, x_i \rangle - y_i)^2$

SET GRADIENT TO 0

$$\frac{\partial L_s}{\partial w} = \frac{2}{m} \sum_{i=1}^{m} [\langle w, x_i \rangle x_i - x_i y_i]$$

$$\frac{\partial L_s}{\partial w} = 0 \implies \sum_{i=1}^{m} \langle w, x_i \rangle x_i = \sum_{i=1}^{m} y_i x_i$$

$A = \sum_{i=1}^{m} x_i x_i^T$

$b = \sum_{i=1}^{m} y_i x_i \qquad Aw = b \implies w = A^{-1} b$

3. Describe how the process can be extended using regularization to avoid large coefficients

We define the regularization function $R(w) = \lambda \|w\|^2$

$\quad \hookrightarrow$ new problem: $\quad w^* = \arg\min_w \frac{1}{m} \sum_{i=1}^{m} \frac{1}{2} (\langle w, x_i \rangle - y_i)^2 + \lambda \|w\|^2$

$$\frac{\partial L_s}{\partial w} = \frac{1}{m} \sum_{i=1}^{m} [\langle w, x_i \rangle x_i - x_i y_i] + 2\lambda w$$

$$\frac{\partial L_s}{\partial w} = 0 \implies 2\lambda m w + \sum_{i=1}^{m} \langle w, x_i \rangle x_i = \sum_{i=1}^{m} y_i x_i$$

$$(2m\lambda I + A) w = b \implies w = (2m\lambda I + A)^{-1} b$$

## CALCULUS

. Define the clustering problem

$X \subseteq \mathbb{R}^d$ : domain

$d: X^2 \to \mathbb{R}^+$ : distance function
→ symmetric : $d(x,y) = d(y,x) \quad \forall x, y \in X$
→ positive
→ triangular inequality : $d(x,z) \leq d(x,y) + d(y,z) \quad \forall x, y, z \in X$

OUTPUT : $C = (C_1, \dots C_k)$

$$\to \bigcup_{i=1}^{k} C_i = X$$

$$\to C_i \cap C_j = \emptyset \quad \forall i, j \in \{1, \dots k\}, i \neq j$$

SOMETIMES $k$ IS GIVEN AS INPUT, SOMETIMES IT IS FREE

2. Introduce the cost function for $k$-means and describe Lloyd's algorithm

$d(x,y) = \|x - y\|_2$ (Euclidean distance)

cluster cost : $\sum_{x \in C_i} d(x, \mu_i)^2$

$\llcorner$ CENTROID

$\mu_i = \dfrac{1}{|C_i|} \sum_{x \in C_i} x$

Lloyd's algorithm :
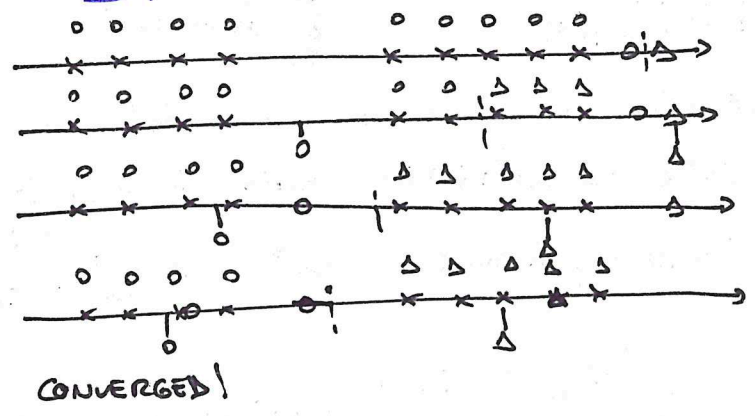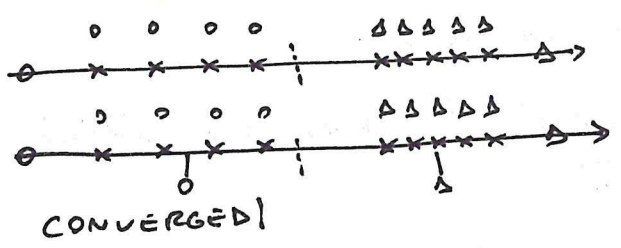1. Select $k$ random centroids
2. Partition the points : $c(x) = \arg\min_{i \in \{1, \dots k\}} d(x, \mu_i) \quad \forall x$
REPEAT $\begin{array}{l} 3. \text{ Recompute the centroids} \end{array}$
4. Convergence: no change / threshold error / number of iterations

3. Use Lloyd's algorithm to solve the clustering problem



CONVERGED!



CONVERGED!

. Describe the classification task

$\quad$ $X$ : domain set

$\quad$ $Y$ : label set

$\quad$ $H : X \to Y$ : hypothesis set

$\quad$ $S$ : training set $\quad ((x_1, y_1), -- (x_m, y_m))$

$\quad$ $f$ : function to be learned (unknown)

$\quad$ $D$ : sampling distribution over $X$

$$L_{D,F} \overset{def}{=} P_{x \sim D} [h(x) \neq f(x)]$$

$$L_S \overset{def}{=} \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}(h(x) \neq f(x))$$

$$\text{ERM:} \quad h^* = \arg\min_{h \in H} L_S(h)$$

2. Describe logistic regression

$$H : \Phi_{sig} \circ L_d$$
$$\quad \overset{\llcorner}{\to} \text{LINEAR}$$
$$\overset{\llcorner}{\to} \text{SIGMOID}$$

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



$$h_w(x) = \frac{1}{1 + e^{-\langle w, x \rangle}}$$
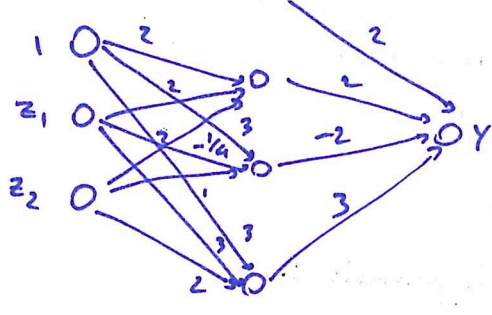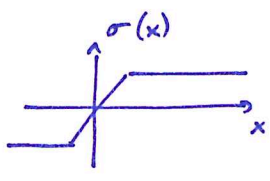$$\quad \overset{\llcorner}{\to} \text{HOMOG. COORD.}$$

$$\ell(h_w, (x, y)) = \log(1 + e^{-y \langle w, x \rangle})$$

$$\text{ERM:} \quad w^* = \arg\min_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^{m} \log(1 + e^{-y_i \langle w, x_i \rangle})$$

$$\sigma(x) = \begin{cases} 1 & x \geq 1 \\ x & -1 \leq x \leq 1 \\ -1 & x \leq -1 \end{cases}$$



$$w^{(1)} = \begin{pmatrix} 2 & 3 & 3 \\ 2 & -\frac{1}{4} & 3 \\ 2 & 1 & 2 \end{pmatrix} \qquad w^{(2)} = \begin{pmatrix} 2 \\ 2 \\ -2 \\ 3 \end{pmatrix}$$

$$z = (1 \quad 3)$$

$$d_{1,1} = \langle w_1^{(1)}, [1, z] \rangle = (2 \; 2 \; 2) \cdot (1 \; 1 \; 3)^T = 7 \qquad o_{1,1} = \sigma(d_{1,1}) = 1$$

$$d_{1,2} = \langle w_2^{(1)}, [1, z] \rangle = (3 \cdot \tfrac{1}{4} \; 1) \cdot (1 \; 1 \; 3)^T = \frac{23}{4} \qquad o_{1,2} = \sigma(d_{1,2}) = 1$$

$$d_{1,3} = \langle w_3^{(1)}, [1, z] \rangle = (3 \; 3 \; 2) \cdot (1 \; 1 \; 3)^T = 12 \qquad o_{1,3} = \sigma(d_{1,3}) = 1$$

$$d_{2,1} = \langle w^{(2)}, [1, o_1] \rangle = (2 \; 2 \; -2 \; 3) \cdot (1 \; 1 \; 1 \; 1)^T = 5 \qquad y = \sigma(d_{2,1}) = 1$$

. Introduce the concept of kernel and its use in SVMs

Feature space transf. : $\psi: X \to F \hookrightarrow$ Hilbert space

$$S = ((x_1, y_1). - (x_-, y_-)) \to \hat{S} ((\psi(x_1), y_1). - (\psi(x_-), y_-))$$

$\mathcal{H}$ : $\underline{\text{linear}}$

Kernel: $K(x, y) = \langle \psi(x), \psi(y) \rangle$

$\hookrightarrow$ can often be computed without going through $\psi$

SVM: $w^* = \overset{\text{arg min}}{\underset{w}{\min}} \left( F(\langle w, \psi(x_1) \rangle \cancel{\}} - \cancel{}\langle w, \psi(x_-) \rangle) + R(\|w\|) \right)$

Representer theor.: $\exists a \in \mathbb{R}^- : w^* = \sum_{i=1}^{-} a_i \psi(x_i)$

$$a^* = \arg\min_{a} \left( F\left( \sum_{j=1}^{-} a_j K(x_j, x_1), - \sum_{j=1}^{-} a_j K(x_j, x_-) \right) + R\left( \sqrt{\sum_{i,j} a_i a_j K(x_i, x_j)} \right) \right)$$

2. Consider the configuration of training data and a scalar function $\Phi(\cdot): \mathbb{R}^2 \to \mathbb{R}$ that makes data linearly separable and relate the map to a kernel



$\Phi(x) = \|x\|_2$
$= x_1^2 + x_2^2$

ONLY POSITIVE

$\langle \Phi(x), \Phi(y) \rangle = x_1^2 y_1^2 + x_1^2 y_2^2 + x_2^2 y_1^2 + x_2^2 y_2^2$

$K_\Phi(x, y) = \langle \langle x, x \rangle, \langle y, y \rangle \rangle$

$(\langle x, y \rangle)^2 = (x_1 y_1 + x_2 y_2)^2 = \left( x_1^2 y_1^2 + x_2^2 y_2^2 + 2 x_1 y_1 x_2 y_2 \right)$

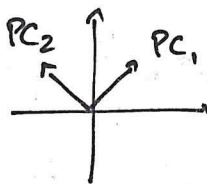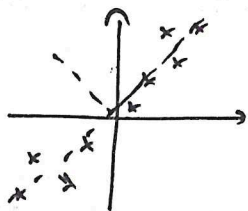. Let $X = [x_1, \ldots, x_n]$, $x_i \in \mathbb{R}^\wedge$ be the data matrix. Introduce PCA

$$W \in \mathbb{R}^{n \times d} \qquad y = Wx$$

$$W = U^T \qquad (U \text{ is an orthonormal base of } \mathbb{R}^n)$$

$$A = \overline{\sum_{i=1}^{n}} x_i x_i^T = VDV^T \qquad (SVD, A \text{ is sym. semidef. pos.})$$

$$PCA: \quad W^* = \left( \underset{U \in \mathbb{R}^{d \times n}, \, \bar{U}U = I}{\arg\max} \; tr\left( U^T A U \right) \right)^T$$

2. Find the 2 PCs for the graph and describe how PCA can be used to simplify linear regression

1. Describe the linear SVM for classification in the case of non-linearly separable data.
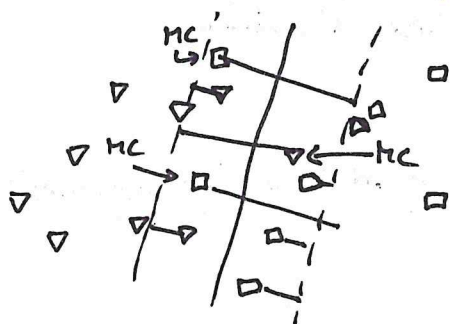
$\xi_i$ : slack variables

$\rightarrow \xi_i \geq 0$

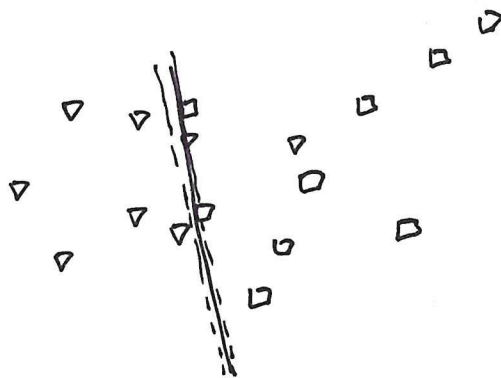$\rightarrow 1 - \xi_i \leq y_i (\langle w, x_i \rangle + b)$

$$\underset{w, b, \xi}{\arg\min} \left( \lambda \|w\|^2 + \frac{1}{m} \sum_{i=1}^{\overline{m}} \xi_i \right)$$

$$\text{s.t. } \xi_i \geq 0, \quad 1 - \xi_i \leq y_i (\langle w, x_i \rangle + b)$$

2. Label the misclassified points and draw the non-zero slack variables. Discuss what happens if $C$ increases.



$C$ increase = $\lambda$ decreases

# Exercise 0

1. Describe the concepts of training and generalization errors in supervised learning

$$\text{TRUE LOSS:} \quad E_{x \sim D}\left[L(h,x)\right]$$

$$\text{TRAINING ERROR:} \quad \frac{1}{L} \sum_{i=1}^{L} L(h, x_i)$$

overfitting: training error is minimized, but true loss grows

2. How would you state the final goal of supervised learning?

→ making predictions that generalize to unseen data: using available data to minimize $L_D$ as well as possible (even if the actual value of $L_D$ is unknowable)

3. What role does k-fold CV play?

Training errors → optimize parameters in a class of hypotheses/algorithm
Validation errors → optimize hyperparameters

k-fold CV (properly implemented) allows to avoid overfitting by having a large hypothesis set.

1. What do you need to show that $VC(H) = d$?

$$C \in X^m$$

$$H_c = \{(h(c_1), \_ h(c_m)), h \in H\}$$

$$h_c \in H_c : C \to \{0,1\}^m$$

$H$ shatters $C$ if $H_c$ contains all possible functions : ~~$H_c = \{0,1\}$~~

$$\exists h_c : C \to \{0,1\}^m \notin H_c$$

$$\boxed{VC(H) = d}$$
1. $\exists C \in X^d$ shattered by $H$
2. $\nexists C \in X^{d+1}$ shattered by $H$ ($\forall C \in X^{d+1}$, $C$ is not shattered)

2. Find the VC dim of right triangles whose legs are parallel to the axes and whose right angle is in the lower left corner

VC dim $\geq 4$ :



(sets of 1 are easy)    (sets of 3 are easy)

the set is shattered

5th point: one point is inside the quadrangle of the other 4



$\rightarrow$ THIS CONFIG IS IMPOSSIBLE!    VCd < 5

$$\boxed{VCd = 4}$$
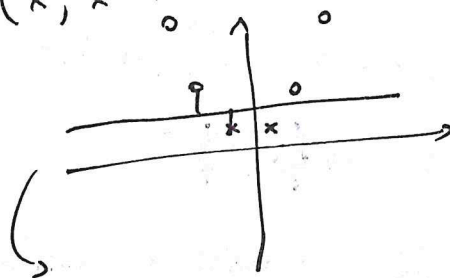
# Exercise 10

| x | y |
|---|---|
| -3 | 1 |
| -2 | 1 |
| -1 | -1 |
| 1 | -1 |
| 2 | 1 |
| 3 | 1 |

$\phi: x \rightarrow (x, x^2)$

$$\text{BOUNDARY}: \quad x^2 = 2.5$$

$$\text{MARGIN}: \quad 1.5$$

| x | $x^2$ | y |
|---|---|---|
| -3 | 9 | 1 |
| -2 | 4 | 1 |
| -1 | 1 | -1 |
| 1 | 1 | -1 |
| 2 | 4 | 1 |
| 3 | 9 | 1 |