# CSC380: Principles of Data Science

**Course wrap-up**

**Xinchen Yu**

- Fill out SCS ([https://scsonline.oia.arizona.edu/](https://scsonline.oia.arizona.edu/)) – if 80% responses, will add 5 points to the homework with lowest grade (currently 55%).

- No lecture next Tuesday, Apr 30
    - You can prepare final exam or work on practice problems in groups and I will do Q&A in person
    - Meinel Optical Sci, Rm 410 (same room)

# Announcements

- ~20 questions and 50% questions will be before midterm.
- Practice questions has been out, keys will be out next week
- No coding questions
- How to prepare
  - **Slides!**
  - Practice problems (helpful but do not only rely on it!)
  - HW questions before midterm

# Probability

Find the Marginal PMFs of X and Y.

$$P_X(0) = \frac{1}{5} + \frac{2}{5} = \frac{3}{5},$$

$$P_X(1) = \frac{2}{5} + 0 = \frac{2}{5},$$

$$P_Y(0) = \frac{1}{5} + \frac{2}{5} = \frac{3}{5},$$

$$P_Y(1) = \frac{2}{5} + 0 = \frac{2}{5}.$$

|  | Y = 0 | Y = 1 |
|---|---|---|
| X = 0 | $\frac{1}{5}$ | $\frac{2}{5}$ |
| X = 1 | $\frac{2}{5}$ | 0 |

# Probability

Find the conditional PMF of X|Y=0 and X|Y=1

$$P_{X|Y}(0|0) = \frac{P_{XY}(0,0)}{P_Y(0)}$$

$$= \frac{\frac{1}{5}}{\frac{3}{5}} = \frac{1}{3}.$$

$$P_{X|Y}(0|1) = 1,$$
$$P_{X|Y}(1|1) = 0.$$

$$P_{X|Y}(1|0) = 1 - \frac{1}{3} = \frac{2}{3}.$$

$$X|Y = 0 \ \sim \ \text{Bernoulli}\left(\frac{2}{3}\right).$$

|  | Y = 0 | Y = 1 |
|---|---|---|
| X = 0 | $\frac{1}{5}$ | $\frac{2}{5}$ |
| X = 1 | $\frac{2}{5}$ | 0 |

# Probability

Let Z=E[X|Y], find the PMF of Z.

|  | Y = 0 | Y = 1 |
|---|---|---|
| X = 0 | $\frac{1}{5}$ | $\frac{2}{5}$ |
| X = 1 | $\frac{2}{5}$ | 0 |

$$Z = E[X|Y] = \begin{cases} E[X|Y = 0] & \text{if } Y = 0 \\ E[X|Y = 1] & \text{if } Y = 1 \end{cases}$$

$$E[X|Y = 0] = \frac{2}{3}, \qquad E[X|Y = 1] = 0,$$

$$Z = E[X|Y] = \begin{cases} \frac{2}{3} & \text{with probability } \frac{3}{5} \\ 0 & \text{with probability } \frac{2}{5} \end{cases}$$

$$P_Z(z) = \begin{cases} \frac{3}{5} & \text{if } z = \frac{2}{3} \\ \frac{2}{5} & \text{if } z = 0 \\ 0 & \text{otherwise} \end{cases}$$

# Probability

Let Z=E[X|Y], find E[Z].

$$P_Z(z) = \begin{cases} \frac{3}{5} & \text{if } z = \frac{2}{3} \\ \frac{2}{5} & \text{if } z = 0 \\ 0 & \text{otherwise} \end{cases}$$

$$E[Z] = \frac{2}{3} \cdot \frac{3}{5} + 0 \cdot \frac{2}{5} = \frac{2}{5}.$$

|  | Y = 0 | Y = 1 |
|---|---|---|
| X = 0 | $\frac{1}{5}$ | $\frac{2}{5}$ |
| X = 1 | $\frac{2}{5}$ | 0 |

# Probability

Let Z=E[X|Y], find var(Z).

$$\text{Var}(Z) = E[Z^2] - (EZ)^2$$
$$= E[Z^2] - \frac{4}{25},$$

$$E[Z^2] = \frac{4}{9} \cdot \frac{3}{5} + 0 \cdot \frac{2}{5} = \frac{4}{15}.$$

$$\text{Var}(Z) = \frac{4}{15} - \frac{4}{25}$$
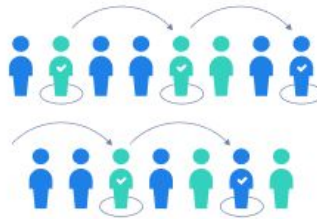$$= \frac{8}{75}.$$

|  | Y = 0 | Y = 1 |
|---|---|---|
| X = 0 | $\frac{1}{5}$ | $\frac{2}{5}$ |
| X = 1 | $\frac{2}{5}$ | 0 |

$$P_Z(z) = \begin{cases} \frac{3}{5} & \text{if } z = \frac{2}{3} \\ \frac{2}{5} & \text{if } z = 0 \\ 0 & \text{otherwise} \end{cases}$$
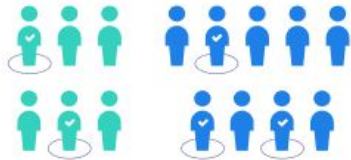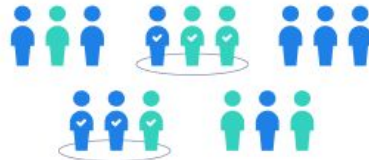
**Simple Random Sample (SRS)**

Each member of the population has the *same chance* of being selected (i.e., uniform over the population)

**Systematic Sample**

Select members of population at a regular interval, determined in advance

**Stratified Sample**

Divide population into *homogeneous* subpopulations (strata). Probability sample the strata.

**Cluster Sample**

Divide population into subgroups (clusters). Randomly select entire clusters.

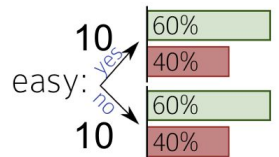# Predictive Modeling and Classification

- Assign all training instances to the root of the tree. Set current node to root node.
- For each feature:
  a. Partition all data instances at the node by the value of the feature.
  b. Compute the accuracy from the partitioning.
- Identify feature that results in the highest accuracy. Set this feature to be the splitting criterion at the current node.
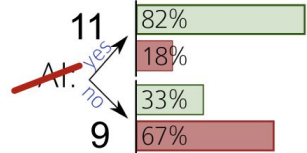
Prereqs  Lecturer  HasLabs

| Rating | Easy? | AI? | Sys? | Thy? | Morning? |
|--------|-------|-----|------|------|----------|
| +2 | y | y | n | y | n |
| +2 | y | y | n | y | n |
| +2 | n | y | n | n | n |
| +2 | n | n | n | y | n |
| +2 | n | y | y | n | y |
| +1 | y | y | n | n | n |
| +1 | y | y | n | y | n |
| +1 | n | y | n | y | n |
| 0 | n | n | n | n | y |
| 0 | y | n | n | y | y |
| 0 | n | y | n | y | n |
| 0 | y | y | y | y | y |
| -1 | y | y | y | n | y |
| -1 | n | n | y | y | n |
| -1 | n | n | y | n | y |
| -1 | y | n | y | n | y |
| -2 | n | n | y | y | n |
| -2 | n | y | y | n | y |
| -2 | y | n | y | n | n |
| -2 | y | n | y | n | y |

overall:
60% like
40% nah

easy:
10 — yes
60%
40%
10 — no
60%
40%

HasTakenPrereqs — AI:
11 — yes
82%
18%
9 — no
33%
67%

SameLecturer — systems:
10 — yes
20%
80%
10 — no
100%
0%

HasLabs — theory:
10 — yes
80%
20%
10 — no
40%
60%
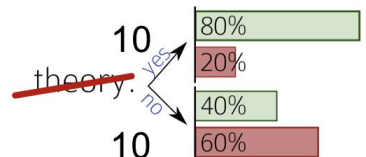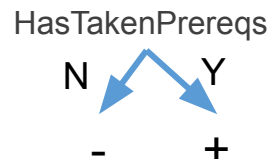
Suppose we place the node HasTakenPrereqs at the root.
Set the prediction at each leaf node as the majority vote.

HasTakenPrereqs

N        Y

-        +

What is the train set accuracy now?

$$\frac{9}{20} \cdot \frac{6}{9} + \frac{11}{20} \cdot \frac{9}{11} = \frac{15}{20} = 0.75$$

No need to split if the leaf is pure
(all data have same labels)

What is the train set accuracy now?

$$\frac{9}{20} \cdot \frac{6}{9} + \boxed{\frac{11}{20}} \cdot \frac{9}{11} = \frac{15}{20} = 0.75$$

Accuracy for two groups:
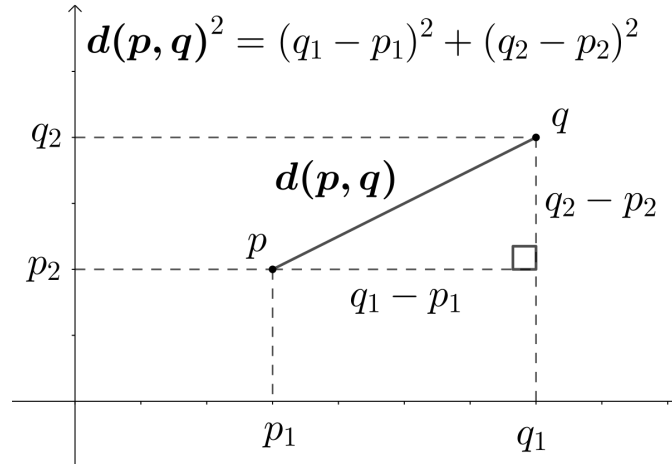- Prereqs = yes (11): 9/11
- Prereqs = no (9): 6/9

For the 11 people prereqs = y, use the majority vote label **like** (9 like, 2 dislike).

Predicted label for 11 people is **like**, 9 people are correctly predicted.

Prereqs   Lecturer   HasLabs

consider it to be 'like'

consider it to be 'dislike'

| Rating | Easy? | AI? | Sys? | Thy? | Morning? |
|---|---|---|---|---|---|
| +2 | y | y | n | y | n |
| +2 | y | y | n | y | n |
| +2 | n | y | n | n | n |
| +2 | n | n | n | y | n |
| +2 | n | y | y | n | y |
| +1 | y | y | n | n | n |
| +1 | y | y | n | y | n |
| +1 | n | y | n | y | n |
| 0 | n | n | n | n | y |
| 0 | y | n | n | y | y |
| 0 | n | y | n | y | n |
| 0 | y | y | y | y | y |
| -1 | y | y | y | n | y |
| -1 | n | n | y | y | n |
| -1 | n | n | y | n | y |
| -1 | y | n | y | n | y |
| -2 | n | n | y | y | n |
| -2 | n | y | y | n | y |
| -2 | y | n | y | n | n |
| -2 | y | n | y | n | y |

# KNN

- Select the number K of the neighbors
- Calculate the Euclidean distance of K number of neighbors
- Take the K nearest neighbors as per the calculated Euclidean distance.
- Among these k neighbors, count the number of the data points in each category.
- Assign the new data points to that category for which the number of the neighbor is maximum.

$$d(p, q)^2 = (q_1 - p_1)^2 + (q_2 - p_2)^2$$

**Training Data:**

| Person | height (feet) | weight (lbs) | foot size(inches) |
|--------|---------------|--------------|-------------------|
| male | 6 | 180 | 12 |
| male | 5.92 (5'11") | 190 | 11 |
| male | 5.58 (5'7") | 170 | 12 |
| male | 5.92 (5'11") | 165 | 10 |
| female | 5 | 100 | 6 |
| female | 5.5 (5'6") | 150 | 8 |
| female | 5.42 (5'5") | 130 | 7 |
| female | 5.75 (5'9") | 150 | 9 |

↑ ↑ ↑

**Features**

**Task:** Observe features $x_1, \ldots, x_D$ and predict class label $y \in \{1, \ldots, C\}$

**Naïve Bayes Model:** Treat features as *conditionally independent* given class label,

$$p(x, y) = p(y)p(x|y) = p(y) \prod_{d=1}^{D} p(x_d \mid y)$$

build individual models for these

**To classify a given instance $x$:** Bayes rule!

$$p(y = c \mid x) = \frac{p(y = c)p(x \mid y = c)}{p(x)}$$

$j : feature, \ c : label, \ i : data$

$y \sim Categorical(\pi_c) : \ p(y = c) = \pi_c$

$$p(y = 1) = \pi_1$$
$$p(y = 2) = \pi_2$$
$$p(y = 3) \ = \pi_3 = 1 - \pi_1 - \pi_2$$

$x|y \sim Bernoulli(\theta_{jc}) : \ p(x|y) = \theta_{jc}{}^x \ (1 - \theta_{jc})^{1-x}$

$x_{j=1}|y = 1 \sim Bernoulli(\theta_{j=1,c=1})$    $x_{j=2}|y = 1 \sim Bernoulli(\theta_{j=2,c=1})$

$x_{j=1}|y = 2 \sim Bernoulli(\theta_{j=1,c=2})$    $x_{j=2}|y = 2 \sim Bernoulli(\theta_{j=2,c=2})$

$x_{j=1}|y = 3 \sim Bernoulli(\theta_{j=1,c=3})$    $x_{j=2}|y = 3 \sim Bernoulli(\theta_{j=2,c=3})$

| $y$ | $x_1$ | $x_2$ |
|---|---|---|
| 1 | 0 | 1 |
| 3 | 1 | 0 |
| 3 | 1 | 1 |
| 2 | 0 | 0 |
| 1 | 1 | 0 |

Q: how many parameters?

c-1+cj

# Model Selection and Evaluation

## **K-fold cross validation**

- Randomly partition train set $S$ into K disjoint sets; call them $\text{fold}_1, \ldots, \text{fold}_K$
- For each hyperparameter $h \in \{1, \ldots, H\}$
    - For each $k \in \{1, \ldots, K\}$
        - train $\hat{f}_k^h$ with $S \setminus \text{fold}_k$
        - measure error rate $e_{h,k}$ of $\hat{f}_k^h$ on $\text{fold}_k$
    - Compute the average error of the above: $\widehat{err}^h = \frac{1}{K}\sum_{k=1}^{K} e_{h,k}$
- Choose $\hat{h} = \arg\min_{h} \widehat{err}^h$
- Train $\widehat{f}^*$ using $S$ (all the training points) with hyperparameter $h$
- Finally, evaluate $\hat{f}^*$ on test set to estimate its future performance.

Use when (1) the dataset is small  (2) ML algorithm's retraining time complexity is low (e.g., kNN)

# 5-fold cross validation



Training Sets      Test Set

Iteration 1    → $Error_1$

Iteration 2    → $Error_2$

Iteration 3    → $Error_3$

Iteration 4    → $Error_4$

Iteration 5    → $Error_5$

Q: If we use 5-fold cross validation for KNN, how many KNNs do we need to train?

A: 5 (if excluding last retraining step).

$$Error = \frac{1}{5}\sum_{i=1}^{5} Error_i$$

**Precision**: dividing the true positives by anything that was predicted as a positive.

$$\frac{\text{TRUE POSITIVES}}{\text{TRUE POSITIVES} + \text{FALSE POSITIVES}}$$

**Recall** (or True Positive Rate): dividing the true positives by anything that should have been predicted as positive.

$$\frac{\text{TRUE POSITIVES}}{\text{TRUE POSITIVES} + \text{FALSE NEGATIVES}}$$

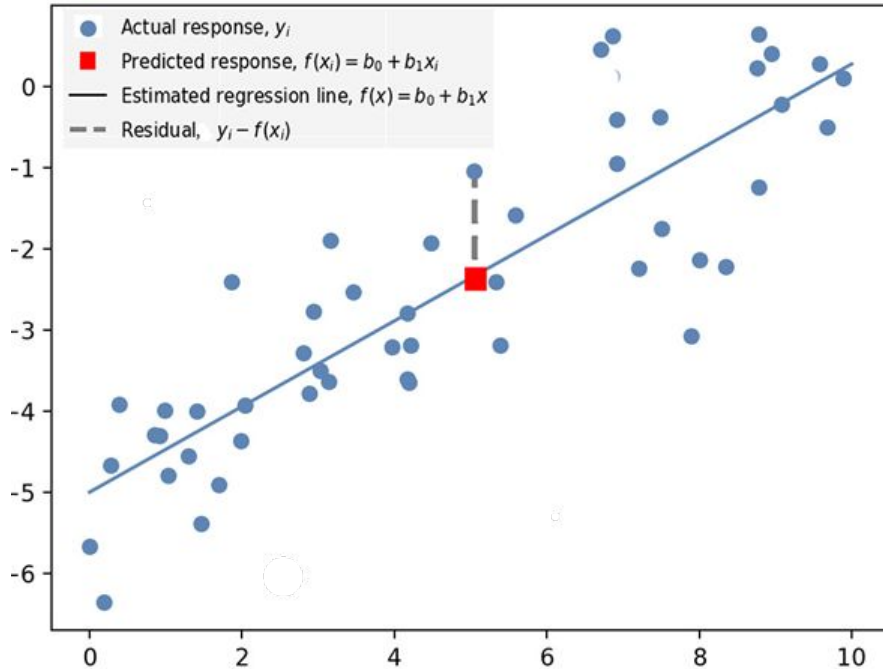F1 score symmetrically represents both precision and recall in one metric.

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{\text{tp}}{\text{tp} + \frac{1}{2}(\text{fp} + \text{fn})}$$

- This is the *harmonic mean* of precision and recall
  - harmonic_mean(x,y)

$$\frac{1}{\frac{1}{2}(\frac{1}{x} + \frac{1}{y})}$$

- Gives equal importance to precision and recall – F1 may not be best when you care about one more than the other (e.g., in medical tests we care about recall)

# Linear Models

Actual response, $y_i$
Predicted response, $f(x_i) = b_0 + b_1 x_i$
Estimated regression line, $f(x) = b_0 + b_1 x$
Residual, $y_i - f(x_i)$

**Functional** Find a line that minimizes the sum of squared residuals!

Given: $\left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1}^{m}$

Compute:

$$w^* = \arg\min_{w} \sum_{i=1}^{m} \left( y^{(i)} - w^T x^{(i)} \right)^2$$

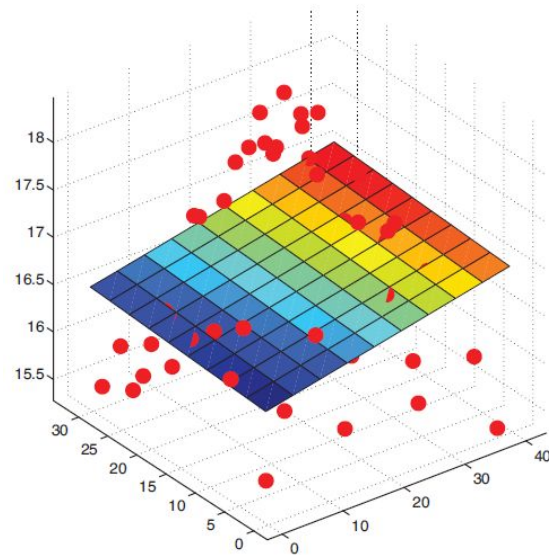*Least squares regression*

Least squares can also be written more compactly,   $\|\boldsymbol{x}\| := \sqrt{\boldsymbol{x} \cdot \boldsymbol{x}}.$

$$\min_{w} \sum_{i=1}^{N} (y^{(i)} - w^T x^{(i)})^2 = \|\mathbf{y} - \mathbf{X}w\|^2$$

Some slightly more advanced linear algebra gives us a solution,

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

***Ordinary Least Squares*** *(OLS)* solution   OLS solution has less residual

# Nonlinear Models

A hyperplane h(x) splits the original
d-dimensional space into two half-spaces.
If the input dataset is linearly separable:

$$y = \begin{cases} +1 & \text{if } h(\mathbf{x}) > 0 \\ -1 & \text{if } h(\mathbf{x}) < 0 \end{cases}$$

Example:

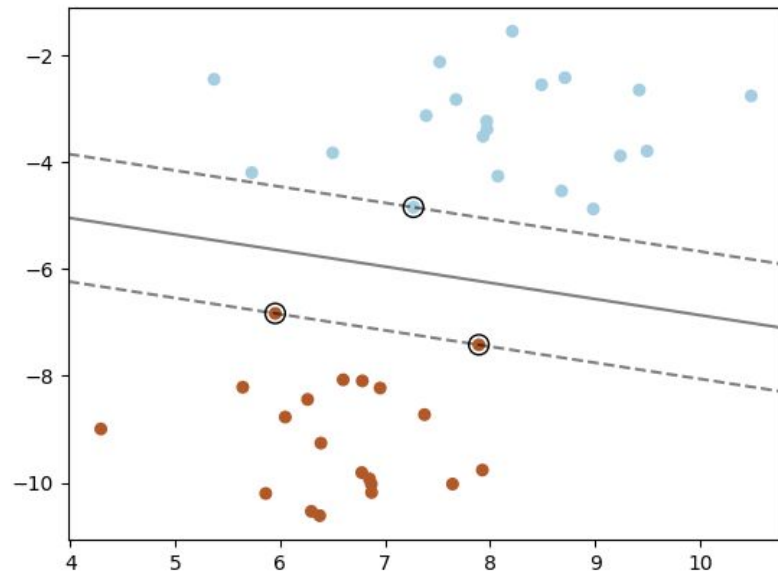$$h(x) = x_1 + 2x_2 - 4$$

Q: label for (0, 3)?

A: +1

Over all the n points, the ***margin*** of the linear classifier is the minimum distance of a point from the separating hyperplane:

$$\delta^* = \min_{\mathbf{x}_i} \left\{ \frac{y_i(\mathbf{w}^T\mathbf{x}_i + b)}{\|\mathbf{w}\|} \right\}$$

All the points that achieve this minimum distance are called ***support vectors***.

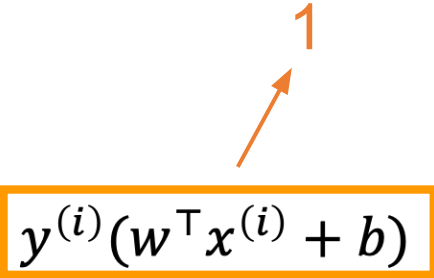$$\delta^* = \frac{y^*(\mathbf{w}^T\mathbf{x}^* + b)}{\|\mathbf{w}\|}$$

Way to solve this issue:
- Choose the scalar s such that the absolute distance of a **support vector** from the hyperplane is 1.

$$sy^*(\mathbf{w}^T\mathbf{x}^* + b) = 1$$

$$s = \frac{1}{y^*(\mathbf{w}^T\mathbf{x}^* + b)} = \frac{1}{y^*h(\mathbf{x}^*)}$$

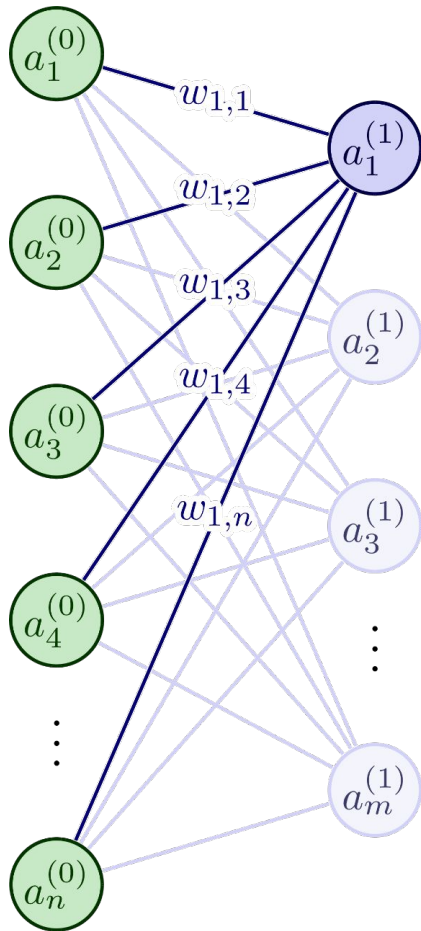$$y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right) \geq 1, \ \text{ for all points } \mathbf{x}_i \in \mathbf{D}$$

$$\underset{w,b}{\arg\max}\ \underset{i}{\min} \boxed{\frac{y^{(i)}(w^\mathsf{T}x^{(i)} + b)}{\|w\|}}$$

1

Q: given a point, how to know if it is support vector?

Margin: $\quad \delta^* = \dfrac{1}{\|\mathbf{w}\|}$

Max margin: $\quad h^* = \arg\max_{h}\{\delta_h^*\} = \arg\max_{\mathbf{w},b}\left\{\dfrac{1}{\|\mathbf{w}\|}\right\}$
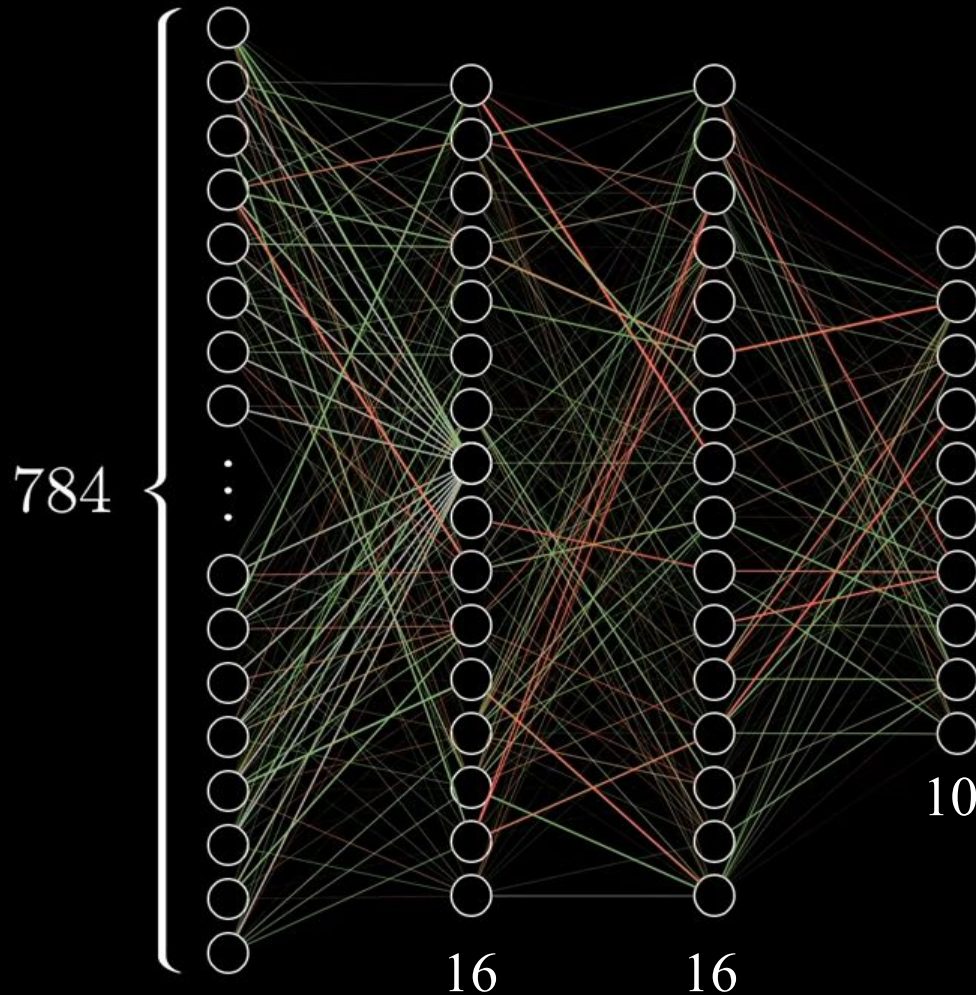
$$= \sigma \left( w_{1,0} a_0^{(0)} + w_{1,1} a_1^{(0)} + \ldots + w_{1,n} a_n^{(0)} + b_1^{(0)} \right)$$

$$= \sigma \left( \sum_{i=1}^{n} w_{1,i} a_i^{(0)} + b_1^{(0)} \right)$$

$$\begin{pmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_m^{(1)} \end{pmatrix} = \sigma \left[ \begin{pmatrix} w_{1,0} & w_{1,1} & \ldots & w_{1,n} \\ w_{2,0} & w_{2,1} & \ldots & w_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m,0} & w_{m,1} & \ldots & w_{m,n} \end{pmatrix} \begin{pmatrix} a_1^{(0)} \\ a_2^{(0)} \\ \vdots \\ a_n^{(0)} \end{pmatrix} + \begin{pmatrix} b_1^{(0)} \\ b_2^{(0)} \\ \vdots \\ b_m^{(0)} \end{pmatrix} \right]$$

$$a^{(1)} = \sigma \left( \mathbf{W}^{(0)} a^{(0)} + \mathbf{b}^{(0)} \right)$$

Number of parameters in this example: $m \times n + m$

$$784 \times 16 + 16 \times 16 + 16 \times 10$$

weights

$$16 + 16 + 10$$

biases

13,002

Each parameter has some impact on the output…need to train all these parameters simultaneously to have a good prediction accuracy
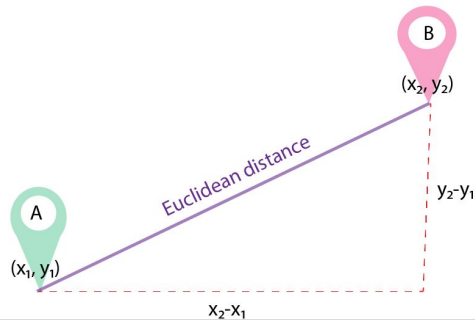
**Input**: $k$: num. of clusters, $S = \{x_1, \dots, x_n\}$

**[Initialize]** Pick $c_1, \dots, c_k$ as randomly selected points from $S$ (see next slides for alternatives)

For t=1,2,…,max_iter

- **[Assignments]**   $\forall x \in S, \quad a_t(x) = \arg\min_{j \in [k]} \|x - c_j\|_2^2$

- If t $\neq 1$  AND  $a_t(x) = a_{t-1}(x), \forall x \in S$
  - break

- **[Centroids]**    $\forall j \in [k], \quad c_j \leftarrow \text{average}(\{x \in S : a_t(x) = j\})$

**Output**: $c_1, \dots, c_k$ and $\{a_t(x_i)\}_{i \in [n]}$

We have the following 3 data points $x_1 = (3, 8)$, $x_2 = (2, 1)$, $x_3 = (5, 4)$. Starting from the initial centroids $c_1 = (0, 0)$ and $c_2 = (4, 3)$, run k-means.

$$d(x_1, c_1)^2 = (3 - 0)^2 + (8 - 0)^2 = 73$$
$$d(x_1, c_2)^2 = (3 - 4)^2 + (8 - 3)^2 = 26$$

$$x_1 \to c_2$$

$$d(x_2, c_1)^2 = 5$$
$$d(x_2, c_2)^2 = 8$$

$$x_2 \to c_1$$

$$d(x_3, c_1)^2 = 41$$
$$d(x_3, c_2)^2 = 2$$

$$x_3 \to c_2$$

Update centroids

$$c_1 = x_2 = (2, 1)$$
$$c_2 = average(x_1, x_3) = \frac{x_1 + x_3}{2} = (4, 6)$$

Stop until centroid remain unchanged