# Machine Learning I
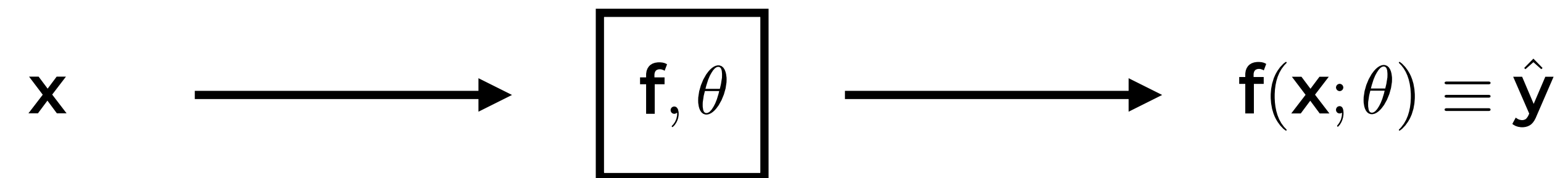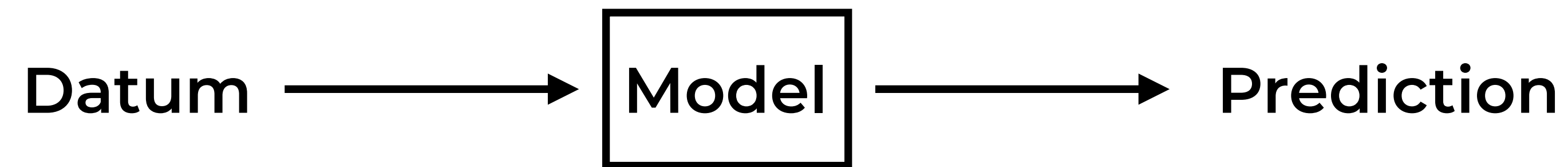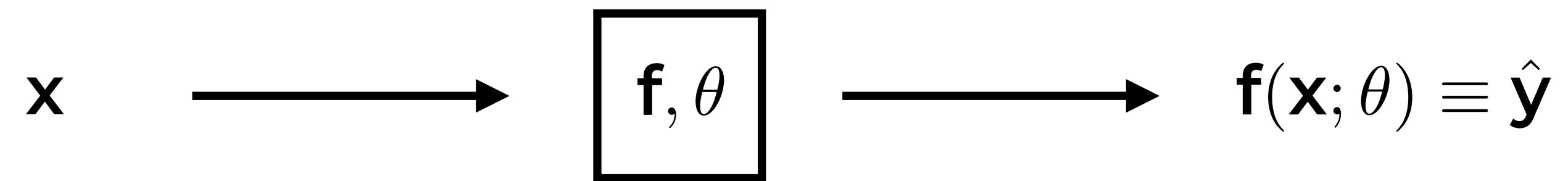# MATH80629A

# Apprentissage Automatique I
# MATH80629

"Mid-term-ish" summary

- **Brief summary of what we have seen so far**

- **Explain concepts within a single framework**

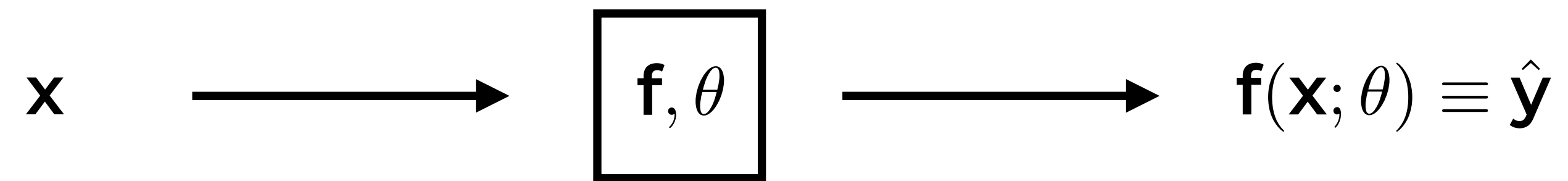- **Focus on a few more advanced concepts**

# Supervised Machine Learning

Datum $\longrightarrow$ Model $\longrightarrow$ Prediction

$\mathbf{x}$ $\longrightarrow$ $\mathbf{f}, \theta$ $\longrightarrow$ $\mathbf{f}(\mathbf{x}; \theta) \equiv \hat{\mathbf{y}}$

# Loss function

$$\mathbf{x} \longrightarrow \boxed{\mathbf{f}, \theta} \longrightarrow \mathbf{f}(\mathbf{x}; \theta) \equiv \hat{\mathbf{y}}$$

# Loss function

Loss

$$\mathbf{x} \longrightarrow \boxed{\mathbf{f}, \theta} \longrightarrow \mathbf{f}(\mathbf{x}; \theta) \equiv \hat{\mathbf{y}}$$

# Loss function

$$\mathsf{x} \quad \longrightarrow \quad \boxed{\mathsf{f}, \theta} \quad \longrightarrow \quad \mathsf{f}(\mathsf{x}; \theta) \equiv \hat{\mathsf{y}}$$

Loss

$$\mathsf{L}(\hat{\mathsf{y}}, \mathsf{y})$$

# Loss function

Loss

$$\mathbf{x} \longrightarrow \boxed{\mathbf{f}, \theta} \longrightarrow \mathbf{f}(\mathbf{x}; \theta) \equiv \hat{\mathbf{y}} \qquad \mathbf{L}(\hat{\mathbf{y}}, \mathbf{y})$$

**Different losses for different types of y's**

| $\mathbf{y} \in \mathcal{R}$ | Regression | $(\hat{\mathbf{y}} - \mathbf{y})^2$ |
|---|---|---|
| **y** categorical e.g., $\{\mathbf{cat}, \mathbf{dog}, \mathbf{bird}\}$ | Classification | accuracy |
| $\mathbf{y} \in \{\mathbf{0}, \mathbf{1}\}$ | Binary Classification | AUC |

# Learning Process

**Distribution over (x,y): P(x,y)**



**Loss**

$x_{train}$ → $\boxed{f, \theta}$ → $\hat{y}_{train}$     $L(\hat{y}_{train}, y_{train})$

$x_{test}$ → $\boxed{f, \hat{\theta}}$ → $\hat{y}_{test}$     $L(\hat{y}_{test}, y_{test})$

# Learning Process
# In practice

**Distribution over (x,y): P(x,y)**



**Loss**

$x_{train} \longrightarrow \boxed{f, \theta} \longrightarrow \hat{y}_{train}$    $L(\hat{y}_{train}, y_{train})$

$x_{valid} \longrightarrow \boxed{f, \hat{\theta}} \longrightarrow \hat{y}_{valid}$

$x_{test} \longrightarrow \boxed{f, \hat{\theta}} \longrightarrow \hat{y}_{test}$    $L(\hat{y}_{test}, y_{test})$

# Learning Process
## In practice

**Distribution over (x,y): P(x,y)**

**Loss**

$x_{train}$ $\longrightarrow$ $\boxed{f, \theta}$ $\longrightarrow$ $\hat{y}_{train}$

$L(\hat{y}_{train}, y_{train})$

$x_{valid}$ $\longrightarrow$ $\boxed{f, \hat{\theta}}$ $\longrightarrow$ $\hat{y}_{valid}$

**Useful:**
**– to select hyper-parameters**
**– To pick the best model**

$x_{test}$ $\longrightarrow$ $\boxed{f, \hat{\theta}}$ $\longrightarrow$ $\hat{y}_{test}$

$L(\hat{y}_{test}, y_{test})$

# Learning

- Learn: Change the parameters to obtain better predictions

Loss

$\theta$

- In other words: change the parameters to minimize the loss

- Take the derivative of the loss wrt the parameter: $\dfrac{d\ Loss}{d\theta}$

# Different models

- f: linear regression, $\theta$ has a closed-form solution

- f: neural network, $\theta$ does not have a closed-forum solution. Gradient descent is used

- Given a training set: $\{(x_{train}, y_{train})\}$

- Initialize $\hat{\theta}_1$ randomly

**for** t $= 1, 2, \ldots$ (epochs) **do**
    **for** i $= 1, 2, \ldots$ (datum) **do**

        - Obtain the predictions $\{f(x_{train}; \hat{\theta}_t)\}$ (Forward propagation)

        - Compute the Loss: $Loss_{ti} := L(f(x_i; \hat{\theta}_t), y_i)$

        - Find the derivative of the loss: $\frac{d\ Loss_{ti}}{d\ \hat{\theta}_t}$

        - Update parameters: $\hat{\theta}_{t+1} = \hat{\theta}_t - \alpha \frac{d\ Loss_{ti}}{d\ \hat{\theta}_t}$

        - If $||\hat{\theta}_{t+1} - \hat{\theta}_t||_2^2 < \epsilon$ then stop
    **end for**
**end for**

- Given a training set: $\{(x_{train}, y_{train})\}$

- Initialize $\hat{\theta}_1$ randomly

**for** $t = 1, 2, \ldots$ (epochs) **do**
    **for** $i = 1, 2, \ldots$ (datum) **do**

Stochastic Gradient Descent

       - Obtain the predictions $\{f(x_{train}; \hat{\theta}_t)\}$ (Forward propagation)

       - Compute the Loss: $\text{Loss}_{ti} := L(f(x_i; \hat{\theta}_t), y_i)$

       - Find the derivative of the loss: $\dfrac{d\ \text{Loss}_{ti}}{d\ \hat{\theta}_t}$
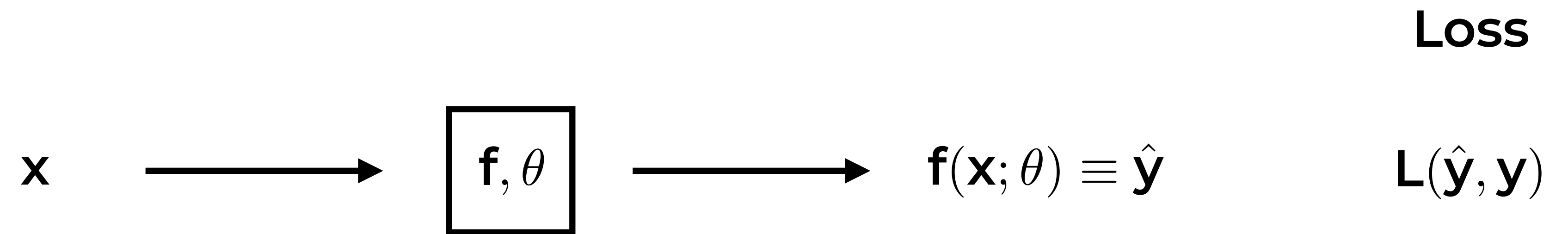
       - Update parameters: $\hat{\theta}_{t+1} = \hat{\theta}_t - \alpha \dfrac{d\ \text{Loss}_{ti}}{d\ \hat{\theta}_t}$

       - If $||\hat{\theta}_{t+1} - \hat{\theta}_t||_2^2 < \epsilon$ then stop
    **end for**
**end for**

# Probabilistic Models

**Loss**

$$\mathbf{x} \longrightarrow \boxed{\mathbf{f}, \theta} \longrightarrow \mathbf{f}(\mathbf{x}; \theta) \equiv \hat{\mathbf{y}} \qquad \mathbf{L}(\hat{\mathbf{y}}, \mathbf{y})$$

# Example

Data: 952  1064  965  1037  871  1029  1138 (unsupervised problem)

# Example

Data: **952  1064  965  1037  871  1029  1138 (unsupervised problem)**

Model: $P(\mathbf{x} \mid \theta) := \mathcal{N}(\boldsymbol{\mu}, \mathbf{1})$

# Example

Data: $\boxed{952}$ 1064 965 1037 871 1029 1138 (unsupervised problem)

Model: $\mathbf{P}(\mathbf{x} \mid \theta) := \mathcal{N}(\boldsymbol{\mu}, \mathbf{1})$

Likelihood for a single datum:

# Example

Data: **952** **1064** **965** **1037** **871** **1029** **1138** (unsupervised problem)

Model: $\mathbf{P}(\mathbf{x} \mid \theta) := \mathcal{N}(\boldsymbol{\mu}, \mathbf{1})$

Likelihood for a single datum:

$$\mathbf{Likelihood}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{1}) = \frac{\mathbf{1}}{\sqrt{\mathbf{2\pi}}} \exp -\frac{(\mathbf{x} - \boldsymbol{\mu})^2}{\mathbf{2}}$$

# Example

Data: 952 1064 965 1037 871 1029 1138 (unsupervised problem)
Model: $P(\mathbf{x} \mid \theta) := \mathcal{N}(\boldsymbol{\mu}, \mathbf{1})$

Likelihood for a single datum:

$$\text{Likelihood}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{1}) = \frac{\mathbf{1}}{\sqrt{\mathbf{2\pi}}} \exp -\frac{(\mathbf{x} - \boldsymbol{\mu})^2}{\mathbf{2}}$$

Log-Likelihood

$$= \log \frac{\mathbf{1}}{\sqrt{\mathbf{2\pi}}} \exp -\frac{(\mathbf{x} - \boldsymbol{\mu})^2}{\mathbf{2}}$$

$$= \mathbf{log1} - \frac{\mathbf{1}}{\mathbf{2}} \mathbf{log2\pi} - \frac{(\mathbf{x} - \boldsymbol{\mu})^2}{\mathbf{2}}$$

# Example

Data: $\boxed{952}$ 1064 965 1037 871 1029 1138 (unsupervised problem)

Model: $P(\mathbf{x} \mid \theta) := \mathcal{N}(\boldsymbol{\mu}, \mathbf{1})$

Likelihood for a single datum:

$$\text{Likelihood}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{1}) = \frac{\mathbf{1}}{\sqrt{\mathbf{2\pi}}} \exp -\frac{(\mathbf{x} - \boldsymbol{\mu})^2}{\mathbf{2}}$$

What value of $\boldsymbol{\mu}$ maximizes it?

**Log-Likelihood**

$$= \log \frac{\mathbf{1}}{\sqrt{\mathbf{2\pi}}} \exp -\frac{(\mathbf{x} - \boldsymbol{\mu})^2}{\mathbf{2}}$$

$$= \mathbf{log1} - \frac{\mathbf{1}}{\mathbf{2}}\mathbf{log2\pi} - \frac{(\mathbf{x} - \boldsymbol{\mu})^2}{\mathbf{2}}$$

# Example

Data: $\boxed{952}$ 1064 965 1037 871 1029 1138 (unsupervised problem)

Model: $P(x \mid \theta) := \mathcal{N}(\mu, 1)$

Likelihood for a single datum:

$\text{Likelihood}(x \mid \mu, 1) = \dfrac{1}{\sqrt{2\pi}} \exp - \dfrac{(x - \mu)^2}{2}$

Log-Likelihood

$= \log \dfrac{1}{\sqrt{2\pi}} \exp - \dfrac{(x - \mu)^2}{2}$

$= \log 1 - \dfrac{1}{2} \log 2\pi - \dfrac{(x - \mu)^2}{2}$

What value of $\mu$ maximizes it?

$\dfrac{d \text{ Log-Likelihood}}{d\,\mu}$

$= \dfrac{d \frac{(x-\mu)^2}{2}}{d\,\mu}$

$= (x - \mu)$

set to 0

$\mu = x$

# Example

Data: $\boxed{952}$ 1064  965  1037  871  1029  1138 (unsupervised problem)

Model: $P(x \mid \theta) := \mathcal{N}(\mu, 1)$

Likelihood for a single datum:

$$\text{Likelihood}(x \mid \mu, 1) = \frac{1}{\sqrt{2\pi}} \exp -\frac{(x - \mu)^2}{2}$$

**Log-Likelihood**

$$= \log \frac{1}{\sqrt{2\pi}} \exp -\frac{(x - \mu)^2}{2}$$

$$= \log 1 - \frac{1}{2} \log 2\pi - \frac{(x - \mu)^2}{2}$$

What value of **μ** maximizes it?

$$\frac{d \text{ Log-Likelihood}}{d \mu}$$

$$= \frac{d \frac{(x-\mu)^2}{2}}{d \mu}$$

$$= (x - \mu)$$

set to 0

$$\mu = x \ = 952$$

**Data: (x,y)**
**Naive Bayes**

# Data: (x,y)
# Naive Bayes

$$\textbf{Model} : \textbf{P}(\textbf{x}, \textbf{y} \mid \theta) = \textbf{P}(\textbf{x} \mid \textbf{y}, \theta)\textbf{P}(\textbf{y} \mid \theta)$$

# Data: (x,y)
# Naive Bayes

$$\mathbf{Model} : \mathbf{P}(\mathbf{x}, \mathbf{y} \mid \theta) = \mathbf{P}(\mathbf{x} \mid \mathbf{y}, \theta)\mathbf{P}(\mathbf{y} \mid \theta)$$

# Data: (x,y)
# Naive Bayes

$$\text{Model}: P(x, y \mid \theta) = P(x \mid y, \theta)P(y \mid \theta)$$

**Data: (x,y)**
**Naive Bayes**

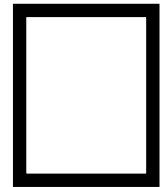$\text{Model}: P(x, y \mid \theta) = P(x \mid y, \theta)P(y \mid \theta)$

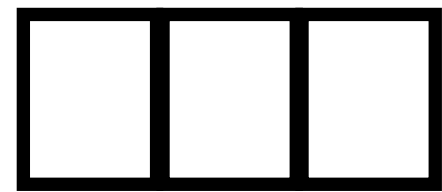**Data: x**
**Gaussian Mixture Models**

# Data: (x,y)
## Naive Bayes

$$\text{Model} : P(x, y \mid \theta) = P(x \mid y, \theta)P(y \mid \theta)$$

# Data: x
## Gaussian Mixture Models

$$\text{Model} : P(x \mid \theta) = \sum_{k=1}^{K} P(\theta_x = k) \underbrace{P(x \mid \theta_k)}_{\mathcal{N}(x \mid \mu_k, \Sigma_k)} \text{ (K components)}$$

# Data: (x,y)
# Naive Bayes

$$\textbf{Model}: \mathbf{P}(\mathbf{x}, \mathbf{y} \mid \theta) = \mathbf{P}(\mathbf{x} \mid \mathbf{y}, \theta)\mathbf{P}(\mathbf{y} \mid \theta)$$

# Data: x
# Gaussian Mixture Models

K=4

$$\textbf{Model}: \mathbf{P}(\mathbf{x} \mid \theta) = \sum_{k=1}^{K} \mathbf{P}(\theta_{\mathbf{x}} = \mathbf{k}) \underbrace{\mathbf{P}(\mathbf{x} \mid \theta_{\mathbf{k}})}_{\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu_k}, \Sigma_{\mathbf{k}})} \textbf{ (K components)}$$
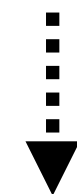
# Data: (x,y)
# Naive Bayes

$$\text{Model} : \mathsf{P}(\mathbf{x}, \mathbf{y} \mid \theta) = \mathsf{P}(\mathbf{x} \mid \mathbf{y}, \theta)\mathsf{P}(\mathbf{y} \mid \theta)$$

# Data: x
# Gaussian Mixture Models

K=4

$$\text{Model} : \mathsf{P}(\mathbf{x} \mid \theta) = \sum_{\mathbf{k=1}}^{\mathbf{K}} \mathsf{P}(\theta_{\mathbf{x}} = \mathbf{k}) \underbrace{\mathsf{P}(\mathbf{x} \mid \theta_{\mathbf{k}})}_{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{\mathbf{k}}, \Sigma_{\mathbf{k}})} \ \text{(K components)}$$

$$\text{Max. likelihood (MLE)} \ : \ \hat{\theta}_{\mathsf{MLE}} = \arg\max_{\theta} \mathsf{P}(\mathbf{x} \mid \theta)$$
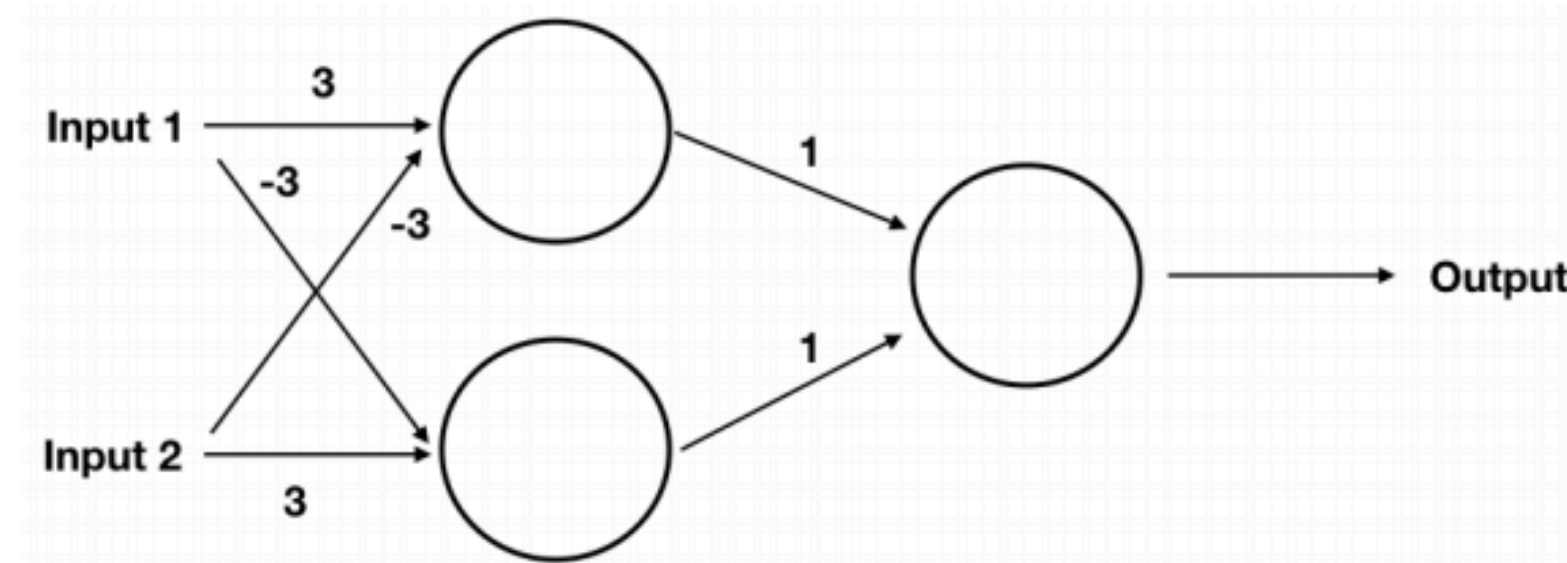
# MLPs / RNNs / CNNs

- MLPs: layers are fully-connected to the next layer

- RNNs: inputs at each layer

  - Typical application: time-series modelling

- CNNs: replace matrix multiplications by convolutions (sparse connections, weight sharing) + pooling

  - Typical application: object recognition in images

# MLPs

(e) (4 points) Consider the neural network below. We have estimated its parameters (shown next to their corresponding arrows).



The activation function of each unit in the network is a simple thresholding function:

$$\text{threshold}(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ 1 & \text{if } x > 0. \end{cases} \tag{1}$$

For each of these four sets of inputs write down the network's output (i.e., its prediction) in the "Output" column of the table below.

| Input 1 | Input 2 | Output |
|---------|---------|--------|
| 0 | 0 | |
| 1 | 1 | |
| 0 | 1 | |
| 1 | 0 | |