# ① Tasks — What do we want to learn? What resources are available?

## Supervised Learning

- Learn to emulate function $x \to y$
- Data tells us correct $y$ for some $x$'s

**Regression**
$y \in \mathbb{R}$

**Classification**
$y$ is discrete

**Binary Classification**
$y \in \{-1, 1\}$

**Multi-class Classification**
$y \in \{1, 2, 3, ..., k\}$

## Unsupervised Learning

- Learn structure of $x$'s
- Dataset only contains bunch of $x$'s

**Embeddings Word Vectors**
For each word, associate it with a vector that capture word similarity, analogy, meaning

**Clustering**
Group structure

**Dimensionality Reduction**
Low dimensional structure

## Reinforcement Learning

- Learn how to take good actions
- Data: After taking action, observe consequences

**Bandits**
State of world is fixed

**Full RL**
Actions change the state of the world

② Modeling — What "shape" does the desired solution have?

- **Tabular Methods**: Remember predicted output for each possible input
  - eg. Tabular Q-Learning
  - Stored a table containing prediction $\hat{Q}(s,a)$ for every state $s$ & action $a$

*Parametric methods*

- **Linear Model**: Can't enumerate all possible inputs
  - So we assume $x \to y$ is a linear function

- **Neural Networks**: Assume $x \to y$ is complex <u>non-linear</u> function
  - MLP: Generic nonlinear function
  - CNN: Local structure matters, weight sharing
  - RNN: Sequential order matters, weight sharing
  - Transformer: Relationships between words matter, so use attention, weight sharing

*add structure compared with basic MLP*

- **Non-parametric Models**: Refer to training data to make predictions
  - K-NN: Similar $x$'s have similar labels
  - <u>Kernel methods</u>: 2 motivations
    - ① Similar $x$'s have similar labels
    - ② Doing linear method in more complex feature space $\phi(x)$

③ **Loss Function**: Quantifies how good/bad a possible solution is

| Supervised Learning | Unsupervised Learning |
|---|---|
| Compare model prediction to true/desired output | want "compressed" version of data make it close to the original |

- **Regression**: $(f(x) - y)^2$
  - model's prediction
  - true answer

- **Binary Classification**: $-\log \sigma(y \cdot f(x))$
  - true label
  - prediction
  - *margin*

- **K-Means**: $\sum_i \| x^{(i)} - \mu_{z_i} \|^2$
  - data
  - assigned cluster mean

- **PCA**: $\sum_i \| x^{(i)} - \text{Proj}_w(x^{(i)}) \|^2$
  - data
  - projection onto subspace spanned by $w$

④ Optimization — How to minimize loss function
                     w.r.t. model's parameters?

Correct → • <u>Direct Computation</u>, e.g. set $\nabla_\theta \text{Loss}(\theta) = 0$
when
possible
                              solve for $\theta$

      — Linear Regression: Normal Equations
      — GMM M-step: Choose best $N, \Sigma, \pi$
                          for each cluster

   • <u>Gradient Descent</u>: $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla_\theta \text{Loss}(\theta^{(t)})$

                                                     learning rate

      — Neural networks
                  (+ backpropagation
                     to compute $\nabla_\theta \text{Loss}(\theta)$ )
      — Linear methods
      — RL: Deep Q Networks, Policy gradient
   • <u>Alternating Minimization</u>:
            — K-Means
            — EM


⑤ Maximum Likelihood Estimation —
          Come up with loss function via probabilistic story
      Loss = negative log likelihood of data
          • Linear / Logistic regression
          • GMM
          • Naive Bayes


⑥ <u>Importance of Data</u> — determines what you learn
          • RL: Choice of action determines what you learn
                think about doing exploration
          • Spurious Correlations
          • Overfitting: too little data → more overfitting
                          ↳ can mitigate w/
                             regularization or simpler model
                  No substitute for more data