

## Algorithm (kNN):

1. Find k examples  $\{\mathbf{x}^{(i)}, t^{(i)}\}$  closest to the test instance  $\mathbf{x}$
2. Classification output is majority class

$$y = \arg \max_{t^{(z)}} \sum_{r=1}^k \delta(t^{(z)}, t^{(r)})$$

**For small k:** overfit. **For large k:** underfit.

**Curse of dimensionality:**  $d$  high  $\rightarrow$  most pts are far away & approx. Same distance (distance variance is very low)

**Normalizations:** nearest neighbors can be sensitive to the ranges of different features.

Without normalization, the performance of kNN will be reduced.

## Decision Tree:

### Entropy:

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x) \quad H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x, y)$$

$$H(Y|X = \text{raining}) = - \sum_{y \in Y} p(y|\text{raining}) \log_2 p(y|\text{raining})$$

$$\begin{aligned} H(Y|X) &= \mathbb{E}_{X \sim p(x)} [H(Y|X)] \\ &= \sum_{x \in X} p(x) H(Y|X = x) \\ &= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(y|x) \\ &= - \mathbb{E}_{(X, Y) \sim p(x, y)} [\log_2 p(Y|X)] \end{aligned}$$

**Chain rule:**  $H(X, Y) = H(X|Y) + H(Y) = H(Y|X) + H(X)$

If X and Y are independent, then X doesn't tell us anything about Y:

$H(Y|X) = H(Y)$ ; But Y tell us everything about Y:  $H(Y|Y) = 0$ .

By knowing X, we can only decrease uncertainty about Y:  $H(Y|X) \leq H(Y)$

$$\begin{aligned} H(X, Y) &= - \sum_{x \in \text{Im}(X)} \sum_{y \in \text{Im}(Y)} p(x, y) \log p(x, y) \\ &= - \sum_{x \in \text{Im}(X)} \sum_{y \in \text{Im}(Y)} p(x, y) \log [p(y) p(x|y)] \\ &= - \sum_{x \in \text{Im}(X)} \sum_{y \in \text{Im}(Y)} p(x, y) \log p(y) - \sum_{x \in \text{Im}(X)} \sum_{y \in \text{Im}(Y)} p(x, y) \log p(x|y) \\ &= - \sum_{y \in \text{Im}(Y)} p(y) \log p(y) - \sum_{x \in \text{Im}(X)} \sum_{y \in \text{Im}(Y)} p(x, y) \log p(x|y) \\ &= H(Y) + H(X|Y). \end{aligned}$$

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

Example:

How much information about cloudiness do we get by discovering whether it is raining?

$$IG(Y|X) = H(Y) - H(Y|X) \approx 0.25 \text{ bits}$$

$$H(Y) = - \left( \frac{49}{100} \log \left( \frac{49}{100} \right) + \frac{51}{100} \log \left( \frac{51}{100} \right) \right)$$

$$H(Y|X) = \sum_{x \in X} p(x) H(Y|X = x) = \frac{25}{100} * \left( - \frac{24}{25} \log \left( \frac{24}{25} \right) - \frac{1}{25} \log \left( \frac{1}{25} \right) \right)$$

$$+ \frac{75}{100} * \left( - \frac{25}{75} \log \left( \frac{25}{75} \right) - \frac{50}{75} \log \left( \frac{50}{75} \right) \right) \quad \text{So } IG(Y|X) = H(Y) - H(Y|X) = 0.25$$

If X is completely uninformative about Y: **IG(Y|X) = 0**

If X is completely informative about Y: **IG(Y|X) = H(Y)** (i.e. in this case,  $H(Y|X) = 0$ )

### Probability rules:

$$\text{Bayes: } P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Correlation  $\rho = \text{Cov}(x, y)$ ,

so  $\text{Cov}(x, y) = \rho \cdot \sigma^2$

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$

$$\text{Var}(aX+b) = a^2 \text{Var}(X)$$

**Decision Tree Construction:** Start with empty decision tree and complete training set. – Split on the most informative attribute, partitioning dataset. – Recurse on sub partitions. **Decision Tree miscellany:** -Exponentially less data at lower levels – Too big of a tree can overfit the data -**Greedy algorithm don't necessarily yield the global optimum** – Mistakes at top-level propagate down tree.

**Decision Tree's advantage over knn:** 1. Good with discrete attributes 2. Easily deals with missing values 3. Robust to scale of inputs, only depends on ordering (kNN has curse of dimensionality, but D.T. doesn't). 4. Fast at test time 5. More interpretable.

**kNN's advantage over decision tree:** 1. Able to handle attributes/features that interact in complex ways (e.g. pixels)

2. Can incorporate interesting distance measures

### Bias-Variance Decomposition:

- Suppose the training set  $\mathcal{D}$  consists of  $N$  pairs  $(\mathbf{x}^{(i)}, t^{(i)})$  sampled independent and identically distributed (i.i.d.) from a sample generating distribution  $p_{\text{sample}}$ ; i.e.,  $(\mathbf{x}^{(i)}, t^{(i)}) \sim p_{\text{sample}}$ .

- Let  $p_{\text{dataset}}$  denote the induced distribution over training sets, i.e.  $\mathcal{D} \sim p_{\text{dataset}}$

- Assume (for the moment) that  $t$  is deterministic given  $\mathbf{x}$ !
- There is a distribution over the loss at  $\mathbf{x}$ , with expectation  $\mathbb{E}_{\mathcal{D} \sim p_{\text{dataset}}} [L(h_{\mathcal{D}}(\mathbf{x}), t)]$ .
- For each query point  $\mathbf{x}$ , the expected loss is different. We are interested in quantifying how well our classifier does over the distribution  $p_{\text{sample}}$ , averaging over training sets:  $\mathbb{E}_{\mathbf{x} \sim p_{\text{sample}}, \mathcal{D} \sim p_{\text{dataset}}} [L(h_{\mathcal{D}}(\mathbf{x}), t)]$ .

$$\begin{aligned} \mathbb{E}_{\mathbf{x}, \mathcal{D}} [(h_{\mathcal{D}}(\mathbf{x}) - t)^2] &= \mathbb{E}_{\mathbf{x}, \mathcal{D}} [(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] + \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] - t)^2] \\ &= \mathbb{E}_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} [(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}])^2 + (\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] - t)^2 + 2(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}]) (\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] - t) | \mathbf{x}]] \\ &= \underbrace{\mathbb{E}_{\mathbf{x}, \mathcal{D}} [(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}])^2]}_{\text{variance}} + \underbrace{\mathbb{E}_{\mathbf{x}} [(\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] - t)^2]}_{\text{bias}} \end{aligned}$$

$$\mathbb{E}_{\mathcal{D}} [2(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}]) (\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] - f(\mathbf{x})) | \mathbf{x}] = 2(\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] - f(\mathbf{x})) \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] | \mathbf{x}].$$

But  $\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] | \mathbf{x}] = 0$ . So:

$$\mathbb{E}_{\mathcal{D}} [2(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}]) (\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) | \mathbf{x}] - f(\mathbf{x})) | \mathbf{x}] = 0.$$

**Bagging:** Suppose we could somehow sample  $m$  independent training sets  $\{\mathcal{D}_i\}_{i=1}^m$  from  $p_{\text{dataset}}$ . We could then learn a predictor  $h_i := h_{\mathcal{D}_i}$  based on each one, and take the average of all the predictors.  $h = \frac{1}{m} \sum_{i=1}^m h_i$

### Bias: Unchanged

$$\mathbb{E}_{\mathcal{D}_1, \dots, \mathcal{D}_m \stackrel{iid}{\sim} p_{\text{dataset}}} [h(\mathbf{x})] = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\mathcal{D}_i \sim p_{\text{dataset}}} [h_i(\mathbf{x})] = \mathbb{E}_{\mathcal{D} \sim p_{\text{dataset}}} [h_{\mathcal{D}}(\mathbf{x})]$$

### Variance: Reduced

$$\text{Var}_{\mathcal{D}_1, \dots, \mathcal{D}_m} [h(\mathbf{x})] = \frac{1}{m^2} \sum_{i=1}^m \text{Var}_{\mathcal{D}_i} [h_i(\mathbf{x})] = \frac{1}{m} \text{Var}_{\mathcal{D}} [h_{\mathcal{D}}(\mathbf{x})].$$

**Problem of bagging:** In reality we don't have  $p_{\text{sample}}$  that can randomly generate datas. Solution: Take a single dataset  $\mathcal{D}$  with  $n$  examples. Then generate  $m$  new datasets, each by sampling  $n$  training examples from  $\mathcal{D}$ , with replacement.

- Problem: the datasets are not independent, so we don't get the  $1/m$  variance reduction.
  - Possible to show that if the sampled predictions have variance  $\sigma^2$  and correlation  $\rho$ , then

$$\text{Var} \left( \frac{1}{m} \sum_{i=1}^m h_i(\mathbf{x}) \right) = \frac{1}{m} (1 - \rho) \sigma^2 + \rho \sigma^2.$$

$\begin{aligned} \text{Var}(y_{\text{avg}}) &= \text{Var} \left( \frac{1}{m} \sum_{i=1}^m y_i \right) \\ &= \frac{1}{m^2} \text{Var} \left( \sum_{i=1}^m y_i \right) \\ &= \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \text{Cov}(y_i, y_j) \\ &= \frac{1}{m^2} [m\sigma^2 + m(m-1)\rho\sigma^2] \\ &= \frac{1}{m} \sigma^2 + \frac{m-1}{m} \rho \sigma^2 \end{aligned}$	<p>Therefore, we use <b>random forest. (bagged decision trees, with when choosing each node, choose a rand set of d input features.</b> It reduces the variance by <b>reducing the correlation between the trees.</b> (i.e. between each model)</p>
--	---

### Bayes Optimality

- Let's return to quantifying expected loss and make the situation slightly more complicated (and realistic): what if  $t$  is not deterministic given  $\mathbf{x}$ ? i.e. have  $p(t|\mathbf{x})$
- We can no longer measure bias as expected distance from true target, since there's a distribution over targets!
- Instead, we'll measure distance from  $y_*(\mathbf{x}) = \mathbb{E}[t | \mathbf{x}]$ 
  - This is the best possible prediction, in the sense that it minimizes the expected loss

Want to show:  $\arg \min_y \mathbb{E}[(y - t)^2 | \mathbf{x}] = y_*(\mathbf{x}) = \mathbb{E}[t | \mathbf{x}]$  (Distribution of  $t \sim p(t|\mathbf{x})$ )

- **Proof:** Start by conditioning on (fixing)  $\mathbf{x}$ .

$$\begin{aligned} \mathbb{E}[(y - t)^2 | \mathbf{x}] &= \mathbb{E}[y^2 - 2yt + t^2 | \mathbf{x}] \\ &= y^2 - 2y\mathbb{E}[t | \mathbf{x}] + \mathbb{E}[t^2 | \mathbf{x}] \\ &= y^2 - 2y\mathbb{E}[t | \mathbf{x}] + \mathbb{E}[t | \mathbf{x}]^2 + \text{Var}[t | \mathbf{x}] \\ &= y^2 - 2yy_*(\mathbf{x}) + y_*(\mathbf{x})^2 + \text{Var}[t | \mathbf{x}] \\ &= (y - y_*(\mathbf{x}))^2 + \text{Var}[t | \mathbf{x}] \end{aligned}$$

- The first term is nonnegative, and can be made 0 by setting  $y = y_*(\mathbf{x})$ .
- The second term doesn't depend on  $y$ ! Corresponds to the inherent unpredictability, or **noise**, of the targets, and is called the **Bayes error** or **irreducible error**.
  - This is the best we can ever hope to do with any learning algorithm. An algorithm that achieves it is **Bayes optimal**.

**Linear Regression:** Optimal weights:  $\mathbf{w}^{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$

**Polynomial curve fitting:** Fit the data using a degree- $M$  polynomial function of the form  $y = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{i=0}^M w_i x^i$

Consider the expression inside the expectation. Since given  $\mathcal{D}$ ,  $h_{\mathcal{D}}$  is determined we can treat the expectation inside the expectation, like 3.1, which gives

$$\mathbb{E}_{\mathcal{D}} [\mathbb{E}_{\mathbf{x}} [y | \mathbf{x}] - y]^2 = \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{\mathbf{x}} [(h_{\mathcal{D}}(\mathbf{x}) - y)^2 | \mathcal{D}]] = \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{\mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})^2 | \mathcal{D}] - 2y \mathbb{E}_{\mathbf{x}} [h_{\mathcal{D}}(\mathbf{x}) | \mathcal{D}] + y^2 | \mathcal{D}]]$$

Since  $\mathbb{E}_{\mathbf{x}} [y | \mathbf{x}] = y$  does not depend on  $\mathcal{D}$  and  $h_{\mathcal{D}}(\mathbf{x}) = \mathbb{E}_{\mathcal{D}} [y | \mathbf{x}]$  does not depend on  $y$ , we can write:

$$\mathbb{E}_{\mathcal{D}, \mathbf{x}} [(h_{\mathcal{D}}(\mathbf{x}) - y)^2] = \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})^2] - 2y \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})] + y^2$$

The last term is Bayes error. In the first term we have  $\mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})^2]$  and it depends on  $\mathbf{x}$  deterministically, so just like 2.1 we can decompose it:

$$\mathbb{E}_{\mathcal{D}, \mathbf{x}} [(h_{\mathcal{D}}(\mathbf{x}) - y)^2] = \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})^2] - 2y \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})] + y^2 = \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})^2] - 2y \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})] + y^2$$

Putting all of these terms together we get:

$$\mathbb{E}_{\mathcal{D}, \mathbf{x}} [(h_{\mathcal{D}}(\mathbf{x}) - y)^2] = \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})^2] - 2y \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})] + y^2 = \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})^2] - 2y \mathbb{E}_{\mathcal{D}, \mathbf{x}} [h_{\mathcal{D}}(\mathbf{x})] + y^2$$

This is called **feature mapping**:  $y = w^T \psi(x)$  where  $\psi(x) = [1, x, x^2, \dots]^T$ . **M too small**: underfit. **M too large**: overfit

**Regularization**: Why regularization? Because we want to reduce overfitting on complex model. Why complex model? Data might support complex models!

$\lambda$  penalizes weight value here.  $\mathcal{J}_{\text{reg}}(\mathbf{w}) = \mathcal{J}(\mathbf{w}) + \lambda \mathcal{R}(\mathbf{w}) = \mathcal{J}(\mathbf{w}) + \frac{\lambda}{2} \sum_j w_j^2$   
If you fit training data poorly,  $J$  is large. If your optimal weights have high values,  $R$  is large.

For least squares problem, we have  $\mathcal{J}(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{t}\|^2$

$$\mathbf{w}_{\lambda}^{\text{Ridge}} = \underset{\mathbf{w}}{\text{argmin}} \mathcal{J}_{\text{reg}}(\mathbf{w}) = \underset{\mathbf{w}}{\text{argmin}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{t}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

$$= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t}$$

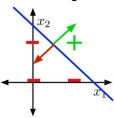
## Logistic Regression:

Binary classification using linear model:

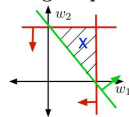
$x_0$	$x_1$	$x_2$	$t$	$z = w_0 x_0 + w_1 x_1 + w_2 x_2$
1	0	0	0	need: $w_0 < 0$
1	0	1	0	need: $w_0 + w_2 < 0$
1	1	0	0	need: $w_0 + w_1 < 0$
1	1	1	1	need: $w_0 + w_1 + w_2 > 0$

AND

Data Space



Weight Space



Verticalization

- Slice for  $x_0 = 1$  and
- example sol:  $w_0 = -1.5, w_1 = 1, w_2 = 1$
- decision boundary:
- $w_0 x_0 + w_1 x_1 + w_2 x_2 = 0$
- $\Rightarrow -1.5 + x_1 + x_2 = 0$

- Slice for  $w_0 = -1.5$  for the constraints
- $w_0 < 0$
- $w_0 + w_2 < 0$
- $w_0 + w_1 < 0$
- $w_0 + w_1 + w_2 > 0$

**Logistic activation function: How do we represent prediction?**

Logistic function is called **sigmoid**.

$$z = \mathbf{w}^T \mathbf{x} + b, \sigma(z) = \frac{1}{1 + e^{-z}}, y = \sigma(z). \text{PS: } \sigma'(z) = \sigma(z) \cdot (1 - \sigma(z))$$

**Loss function for logistic regression:**

$$L_{CE} = -t \log y - (1 - t) \log(1 - y)$$

**Gradient of logistic loss:**

$$\frac{\partial L_{CE}}{\partial w_j} = \frac{\partial L_{CE}}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w_j} = \left( -\frac{t}{y} + \frac{1-t}{1-y} \right) \cdot y(1-y) \cdot x_j = (y - t)x_j$$

Problem: if  $t = 1$  and I am really confident that it's a negative example. i.e.  $z \ll 0$ . Then  $\log(y)$  might be numerically 0.

Solution:  $L_{LCE}(z, t) = L_{CE}(\sigma(z), t) = t \log(1 + e^{-z}) + (1 - t) \log(1 + e^z)$

**Neural Networks:**

**Multiclass classification:** Targets form a discrete set  $\{1, \dots, K\}$ ,  $K$  classes.

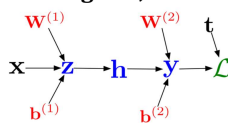
Targets are represented as **1-hot** vectors.  $\mathbf{t} = [0, 0, \dots, 1, 0, \dots, 0]$

Prediction for the probability of each class  $y_k$  is between  $[0, 1]$

$$y_k = \text{softmax}(z_1, \dots, z_K)_k = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}, \text{ where } z_k \text{ is called logits}$$

**CE for multiclass classification:**  $L_{CE}(\mathbf{y}, \mathbf{t}) = -\mathbf{t}^T (\log \mathbf{y})$

**NN's diagram, in both scalar and vector form:**



**Backward pass:**

**Forward pass:**

$$\mathbf{z} = \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}$$

$$z_i = \sum_j w_{ij}^{(1)} x_j + b_i^{(1)}$$

$$\mathbf{h} = \sigma(\mathbf{z})$$

$$h_i = \sigma(z_i)$$

$$\mathbf{y} = \mathbf{W}^{(2)} \mathbf{h} + \mathbf{b}^{(2)}$$

$$y_k = \sum_i w_{ki}^{(2)} h_i + b_k^{(2)}$$

$$\mathcal{L} = \frac{1}{2} \|\mathbf{t} - \mathbf{y}\|^2 = \frac{1}{2} \sum_k (y_k - t_k)^2$$

$$\bar{z} = 1, \quad \bar{y}_k = \bar{z}(y_k - t_k)$$

$$\bar{w}_{ki}^{(2)} = \bar{y}_k h_i, \quad \bar{\mathbf{W}}^{(2)} = \bar{\mathbf{y}} \mathbf{h}^T$$

$$\bar{b}_k^{(2)} = \bar{y}_k, \quad \bar{\mathbf{b}}^{(2)} = \bar{\mathbf{y}}$$

$$\bar{h}_i = \sum_k \bar{y}_k w_{ki}^{(2)}, \quad \bar{\mathbf{h}} = \mathbf{W}^{(2)T} \bar{\mathbf{y}}$$

$$\bar{z}_i = \bar{h}_i \sigma'(z_i), \quad \bar{\mathbf{z}} = \bar{\mathbf{h}} \circ \sigma'(\mathbf{z})$$

$$\bar{w}_{ij}^{(1)} = \bar{z}_i x_j, \quad \bar{\mathbf{W}}^{(1)} = \bar{\mathbf{z}} \mathbf{x}^T$$

$$\bar{b}_i^{(1)} = \bar{z}_i, \quad \bar{\mathbf{b}}^{(1)} = \bar{\mathbf{z}}$$

**GD Rules:** let  $g$  be gradient. Then

GD:  $\mathbf{w} \leftarrow \mathbf{w} - \eta \sum_{i=1}^N g_i$

SGD: Choose  $i \sim U[1, N]$ ,  $\mathbf{w} \leftarrow \eta g_i$

Mini-batch GD: Choose a subset  $M \subset \{1, \dots, N\}$ ,  $\mathbf{w} \leftarrow \eta \sum_{i \in M} g_i$

Computation cost:  $\text{SGD} < \text{mSGD} < \text{GD}$

Convergence rate:  $\text{SGD} < \text{mSGD} < \text{GD}$

**SGD's mathematical justification:** if you sample a training example **uniformly at random**, the stochastic gradient is an **unbiased estimate** of the batch

$$\text{gradient. } E \left[ \frac{\partial L^{(i)}}{\partial \theta} \right] = \frac{1}{N} \sum_{i=1}^N \frac{\partial L^{(i)}}{\partial \theta} = \frac{\partial J}{\partial \theta}$$

## SVM:

Goal: Minimize the hinge loss that leads to an optimal separating hyperplane.

What is the optimal separating hyperplane? A **hyperplane** is described by points  $\mathbf{x} \in \mathbb{R}^D$  s.t.  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$ .

We want to find the hyperplane that minimize the training loss:

$$\min_{\mathbf{w}, b} \sum_{i=1}^N \max\{0, 1 - t^{(i)} z^{(i)}(\mathbf{w}, b)\} \text{ (minimizing the hinge loss! } \max\{0, 1 -$$

$t^{(i)} z^{(i)}(\mathbf{w}, b)\}$  is the hinge loss). Often add  $\frac{\lambda}{2} \|\mathbf{w}\|^2$  in the end (regularizer). **Decision boundary** is linear. Can only do binary classification.

To find  $\mathbf{w}$  and  $b$ : Find all the three points that are support vectors. Then, substitute their coordinates to  $\mathbf{w}^T \mathbf{x} + b$ . The right side equals whether it is a positive example (1) or negative example (-1). E.g.  $1w_1 + 4w_2 + b = 1, 2w_1 + 3w_2 + b = 1, 4w_1 + 5w_2 + b = -1$ . Then solve for this system of linear equation! You will get  $\mathbf{w}$  and  $b$ .

**What is the benefit of the hinge loss compared to 0-1 loss?** Minimum of a function will be at its critical points. But the gradient of the 0-1 loss is zero everywhere it's defines. This means that changing the weights by a very small amount probably has no effect on the loss. We can use the 0-1 loss function, but minimizing this loss is computationally difficult, and it cannot distinguish different hypotheses that achieve the same accuracy. The hinge loss is a relaxation of the 0-1 loss, it can be minimized using gradient descent.

**Claim: hinge loss = 0 iff all points are correctly classified.**

## Boosting:

AdaBoost Algorithm:

- Input: Data  $\mathcal{D}_N$ , weak classifier WeakLearn (a classification procedure that returns a classifier  $h$ , e.g. best decision stump, from a set of classifiers  $\mathcal{H}$ , e.g. all possible decision stumps), number of iterations  $T$
- Output: Classifier  $H(x)$
- Initialize sample weights:  $w^{(n)} = \frac{1}{N}$  for  $n = 1, \dots, N$
- For  $t = 1, \dots, T$ 
  - Fit a classifier to data using weighted samples ( $h_t \leftarrow \text{WeakLearn}(\mathcal{D}_N, \mathbf{w})$ ), e.g.,

$$h_t \leftarrow \underset{h \in \mathcal{H}}{\text{argmin}} \sum_{n=1}^N w^{(n)} \mathbb{I}\{h(\mathbf{x}^{(n)}) \neq t^{(n)}\}$$

- Compute weighted error  $\text{err}_t = \frac{\sum_{n=1}^N w^{(n)} \mathbb{I}\{h_t(\mathbf{x}^{(n)}) \neq t^{(n)}\}}{\sum_{n=1}^N w^{(n)}}$
- Compute classifier coefficient  $\alpha_t = \frac{1}{2} \log \frac{1 - \text{err}_t}{\text{err}_t} \in (0, \infty)$
- Update data weights

$$w^{(n)} \leftarrow w^{(n)} \exp \left( -\alpha_t t^{(n)} h_t(\mathbf{x}^{(n)}) \right) \left[ \equiv w^{(n)} \exp \left( 2\alpha_t \mathbb{I}\{h_t(\mathbf{x}^{(n)}) \neq t^{(n)}\} \right) \right]$$

- Return  $H(\mathbf{x}) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right)$

Weak classifiers which get lower weighted error ( $\text{err}_t$ ) get more weight in the final classifier. (Since  $\text{err}_t$  low  $\Rightarrow \alpha_t$  high, so when computing  $H(\mathbf{x})$  it gets more weight) Also, when updating weights:  $w^{(n)} \leftarrow w^{(n)} \exp(2\alpha_t \mathbb{I}\{h_t(\mathbf{x}^{(n)}) \neq t^{(n)}\})$

If  $\text{err}_t \approx 0$ ,  $\alpha_t$  high so misclassified examples get more attention

If  $\text{err}_t \approx 0.5$ ,  $\alpha_t$  low so misclassified examples get less attention

(i.e. stronger classifiers misclassified  $\Rightarrow$  add more weight!)

AdaBoost Minimizes the training error.

**AdaBoost penalizes more regarding error compared to logistic regression.**

**Theorem:** Assume that at each iteration of AdaBoost the WeakLearn returns a hypothesis with error  $\text{err}_t \leq \frac{1}{2} - \gamma$  for all  $t = 1, \dots, T$  with  $\gamma > 0$ . The training error of the output hypothesis  $H(\mathbf{x}) = \text{sign}(\sum_{t=1}^T \alpha_t h_t(\mathbf{x}))$  is at most

$L_N(H) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{H(\mathbf{x}^{(i)}) \neq t^{(i)}\} \leq \exp(-\gamma^2 T)$ . This is under the simplifying assumption that each weak learner is  $\gamma$ -better than a random predictor. It is called a geometric convergence.  $\rightarrow$  Fast! **AdaBoost is an additive model.**

**Q: How does AdaBoost reduce bias?** By making each classifier focus on previous mistakes **Q: Sometimes the test error decreases even after the training error is zero. How does that happen?** The confidence of the classifier is not just the training error. After the training error is zero, the confidence continues to increase as each example goes further away from the decision boundary. The farther an example is, the more confident the model is in its label, or another way of thinking about it, the decision is more robust to noise in the input features.

**Additive models:** We obtain the additive model  $H_m(\mathbf{x}) = \sum_{i=1}^m \alpha_i h_i(\mathbf{x})$  with

$$h_m \leftarrow \underset{h \in \mathcal{H}}{\text{argmin}} \sum_{i=1}^N w_i^{(m)} \mathbb{I}\{h(\mathbf{x}^{(i)}) \neq t^{(i)}\}, \quad \alpha = \frac{1}{2} \log \left( \frac{1 - \text{err}_m}{\text{err}_m} \right),$$

$$\text{where } \text{err}_m = \frac{\sum_{i=1}^N w_i^{(m)} \mathbb{I}\{h_m(\mathbf{x}^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w_i^{(m)}}, \quad w_i^{(m+1)} = w_i^{(m)} \exp \left( -\alpha_m h_m(\mathbf{x}^{(i)}) t^{(i)} \right).$$

**Bagging vs Boosting:**

**Bagging:** Train classifiers independently on random subsamples of training data. 1. reduces variance 2. bias is not changed much 3. parallel 4. want to minimize



correlation between ensemble elements **Boosting:** Train classifier sequentially, each time focusing on training data points that were previously misclassified. 1. reduces bias 2. increases variance 3. sequential 4. high dependency between ensemble elements

**Generative model vs Discriminative model:**

**Discriminative approach:** estimate parameters of decision boundary/class separator directly from labeled examples. **1.** Tries to solve: How do I separate the classes? **2.** Learn  $p(t|\mathbf{x})$  directly (logistic regression models) **3.** Learn mappings from inputs to classes (linear/logistic regression, decision trees etc)

**Generative approach:** model the distribution of inputs characteristic of the class (Bayes classifier) **1.** Tries to solve: What does each class “look” like? **2.** Build a model of  $p(\mathbf{x}|t)$  **3.** Apply Bayes rule i.e.  $p(t, \mathbf{x}) = \dots$

**Key difference:** Is there a distributional assumption over inputs?

**Bayes Classifier:** Given features  $\mathbf{x} = [x_1, x_2, \dots, x_D]^T$  we want to compute class probabilities using Bayes Rule:

$$p(c|\mathbf{x}) = \frac{p(\mathbf{x}, c)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|c)p(c)}{p(\mathbf{x})}$$

More formally:

$$posterior = \frac{Class\ likelihood \times prior}{Evidence}$$

To compute  $p(c|\mathbf{x})$  we need:  $p(\mathbf{x}|c)$  and  $p(c)$

**Naïve Bayes** assumption: Naïve Bayes assumes that the world features  $x_i$  are **conditionally independent** given the class  $c$ .

$$p(c, x_1, \dots, x_D) = p(c)p(x_1|c) \dots p(x_D|c)$$

The log-likelihood can be decomposed into:

$$l(\theta) = \sum_{i=1}^N \log p(c^{(i)}, \mathbf{x}^{(i)}) = \sum_{i=1}^N \log p(c^{(i)}) + \sum_{j=1}^D \sum_{i=1}^N \log p(x_j^{(i)}|c^{(i)})$$

**To do learning:** (during training) Find MLE for  $\pi = p(c^{(i)} = 1)$  and  $\theta = p(x_j^{(i)} = 1|c)$  (Since good values of  $\theta$  should assign high probability to the observed data.)

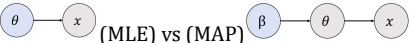
For the prior we maximize:  $\sum_{i=1}^N \log p(c^{(i)})$ . Then  $p(c^{(i)}) = \pi^{c^{(i)}} (1 - \pi^{1-c^{(i)}})$

For the class likelihood  $\theta$  we maximize:  $\sum_{i=1}^N \log p(x_j^{(i)}|c^{(i)})$ . Then  $p(x_j^{(i)}|c = \theta_{jc}^{x_j^{(i)}} (1 - \theta_{jc})^{1-x_j^{(i)}}$

**To do prediction:** (during test time) We predict the category by performing inference in the model. Apply **Bayes’ Rule**:

$$p(c|\mathbf{x}) = \frac{p(c)p(\mathbf{x}|c)}{\sum_{c'} p(c')p(\mathbf{x}|c')} = \frac{p(c) \prod_{j=1}^D p(x_j|c)}{\sum_{c'} p(c') \prod_{j=1}^D p(x_j|c')} \propto p(c) \prod_{j=1}^D p(x_j|c)$$

**MLE’s disadvantage:** Not good with sparse data! e.g., What if you flip the coin **only twice** (so data is sparse) and get H both times? Then  $\theta_{ML} = \frac{N_H}{N_H + N_T} = \frac{2}{2+0} = 1$ , i.e. the probability of getting head is 1, according to the model’s belief! Therefore, it assigned getting tail’s outcome probability 0. : ( Solution: Add prior distribution to theta!)



**MAP Estimation:**  $\theta_{MAP} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \max_{\theta} p(\theta, \mathcal{D}) = \arg \max_{\theta} p(\theta)p(\mathcal{D}|\theta) = \arg \max_{\theta} \log p(\theta) + \log p(\mathcal{D}|\theta)$

**Gaussian Discriminative Analysis (GDA):**

Assumes that  $p(\mathbf{x}|t)$  is distributed according to multivariate Gaussian distribution.

$$p(\mathbf{x}|t = k) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_k|^{\frac{1}{2}}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k) \right]$$

Where  $|\Sigma_k|$  denotes the determinant of the matrix, and  $D$  is dimension of  $\mathbf{x}$ . Each class  $k$  has a mean vector  $\mu_k$  and a covariance matrix  $\Sigma_k$   
Use MLE to find empirical mean and empirical covariance matrix:  $\hat{\mu}_k$  and  $\hat{\Sigma}_k$ .  $\Sigma_k$  has  $O(D^2)$  parameters – could be hard to estimate  
To find log-likelihood of GDA: Assume data  $\mathbf{x}$  are drawn from gaussian, with  $\mu$  and  $\Sigma$ . We calculate log-likelihood  $l(\mu, \Sigma)$ , then take derivative w.r.t  $\mu$  and  $\Sigma$ , evaluate at zero to find the MLE  $\hat{\mu}$  and  $\hat{\Sigma}$ .

**Decision boundary of GDA: usually quadratic.** It is based on class posterior. GDA makes decisions by comparing class probabilities.

$$\log p(t_k|\mathbf{x}) = \log p(\mathbf{x}|t_k) + \log p(t_k) - \log p(\mathbf{x}) \text{ \#By bayes rule}$$
$$= -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_k^{-1}| - \frac{1}{2} (\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k) + \log p(t_k) - \log p(\mathbf{x})$$

Now, suppose I have two classes,  $l$  and  $k$ . Then the decision boundary equals the follow: decision boundary is  $\log p(t_k|\mathbf{x}) = \log p(t_l|\mathbf{x})$ , which is equivalent to

$$(\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k) = (\mathbf{x} - \mu_l)^T \Sigma_l^{-1} (\mathbf{x} - \mu_l) + C_{k,l}$$

Quadratic function in  $\mathbf{x} \Rightarrow$  **quadratic decision boundary!**

What is  $C_{k,l}$ ? What if  $\Sigma_k = \Sigma_l$ ?  $\rightarrow$  **can be derived. Easy.**

**Simplifying GDA model:** What if  $\mathbf{x}$  is high-dimensional?

For Gaussian Bayes Classifier, if input  $\mathbf{x}$  is high-dimensional, then covariance matrix has many parameters  $O(D^2)$ . We can save some parameters by using a shared covariance for the classes, i.e.  $\Sigma_k = \Sigma_l$ . **If so, we have a linear decision boundary!**

**Comparing Naïve Bayes vs GDA:**

Both Naïve Bayes and GDA are generative models. But:

Naïve Bayes	GDA
Works on a set of discrete inputs	Works on continuous inputs

Assume conditional independence given class, so $p(x_1, x_2 c) = p(x_1 c)p(x_2 c)$ . (i.e. Covariance is <b>0!</b> ) Objective is to figure out, for each $x_i$ and each $c$ , the $p(x_i c)$ that maximizes the likelihood of the observed data	Assume the underlying data is generated from Gaussian distribution. This let us specify covariance through covariance matrices, instead of being 0. Also, see the following sentence:
---	--

GDA’s objective is then to figure out the Gaussian parameters (mean and covariance matrix) that provide the best explanation (highest likelihood) for the observed data.

**Comparing GDA vs Logistic Regression:** GDA makes stronger modeling assumption: assumes class-conditional data is multivariate Gaussian. If this is true, GDA is asymptotically efficient. But LR is more robust, less sensitive to incorrect modeling assumptions. When the class-conditional distributions are non-Gaussian (true almost always), LR usually beats GDA. Also, GDA has quadratic (when not shared covariance), LR has linear decision boundary.

**PCA:**

PCA is a linear model. Goal: Dimensionality reduction  $\rightarrow$  saves computation! Need to choose  $D \times K$  matrix  **$U$  with orthonormal columns.**

Two criteria: **1.** Minimize the **reconstruction error:** Find vectors in a subspace that are closest to data points.  $\min_U \frac{1}{N} \sum_{i=1}^N ||\mathbf{x}^{(i)} - \tilde{\mathbf{x}}^{(i)}||^2$  **2.** Maximize the variance of reconstructions: Find a subspace where data has the most variability.

$\max_U \frac{1}{N} \sum_{i=1}^N ||\tilde{\mathbf{x}}^{(i)} - \hat{\mu}||^2$ . **Proof the data and its reconstruction has the same means:**

**To find the best  $\tilde{\mathbf{x}}$  (i.e. to find the best projection):**

Let  $U$  be a matrix with columns  $\{\mathbf{u}_k\}_{k=1}^K$ , then  $\mathbf{z} = U^T (\mathbf{x} - \hat{\mu})$ . Also,  $\tilde{\mathbf{x}} = \hat{\mu} + U\mathbf{z} = \hat{\mu} + UU^T (\mathbf{x} - \hat{\mu})$ . Here,  $UU^T$  is the projector on subspace, and  $U^T U = I$  ( $UU^T \neq I!!!$ )

**Proof that dimensions of  $\mathbf{z}$  are decorrelated:**  $Cov(\mathbf{z}) = Cov(U^T (\mathbf{x} - \hat{\mu})) = U^T Cov(\mathbf{x}) U = U^T \Sigma U = U^T Q \Lambda Q^T U = (I \quad \mathbf{0}) \begin{pmatrix} I \\ \mathbf{0} \end{pmatrix}$  = *top left*  $K \times K$  block of  $\Lambda$ . Since the covariance matrix is diagonal, this means the features are uncorrelated. **PCA can be viewed as a linear autoencoder.** (as shown on the left). If  $K < D$ , then  **$W1$**  maps  $\mathbf{x}$  to a  $K$ -dimensional space, so it’s doing dimensionality reduction.

**Unsupervised Learning:**

**K-means algorithm:** Objective: Find cluster centers  $\{\mathbf{m}_k\}_{k=1}^K$  and assignments  $\{r^{(n)}\}_{n=1}^N$  to minimize the sum of squared distances of data points  $\{\mathbf{x}^{(n)}\}$  to their assigned cluster centers. **1.** Data sample  $n = 1, \dots, N$ :  $\mathbf{x}^{(n)} \in \mathbb{R}^D$  (observed) **2.** Cluster center  $k = 1, \dots, K$ :  $\mathbf{m}_k \in \mathbb{R}^D$  (not observed) **3.** Responsibilities: Cluster assignment for sample  $n$ :  $r^{(n)} \in \mathbb{R}^K$  1-of- $K$  encoding (not observed)  
Mathematically, we want:

$$\min_{\{\mathbf{m}_k\}, \{r^{(n)}\}} J(\{\mathbf{m}_k\}, \{r^{(n)}\}) = \min_{\{\mathbf{m}_k\}, \{r^{(n)}\}} \sum_{n=1}^N \sum_{k=1}^K ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2$$

The algorithm: Initialization: Set  $K$  cluster means  $\mathbf{m}_1, \dots, \mathbf{m}_K$  to random values. Repeat the following until convergence (i.e. until assignments do not change): **Assignment:** Optimize  $J$  w.r.t.  $\{r\}$ : Each data point  $\mathbf{x}^{(n)}$  assigned to nearest center

$$\hat{k}^{(n)} = \arg \min_k ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2$$

And **Responsibilities** (1 hot, if not using soft K-means):

$$r_k^{(n)} = I(\hat{k}^{(n)} = k) \text{ for } k = 1, \dots, K$$

**Refitting:** Optimize  $J$  w.r.t.  $\{\mathbf{m}\}$ : Each center is set to mean of data assigned to it.

$$\mathbf{m}_k = \frac{\sum_n r_k^{(n)} \mathbf{x}^{(n)}}{\sum_n r_k^{(n)}}$$

K-means algorithm reduces the cost **at each iteration.**

**Soft K-means:** different assignment step.

$$r_k^{(n)} = \frac{\exp[-\beta ||\mathbf{m}_j - \mathbf{x}^{(n)}||^2]}{\sum_j \exp[-\beta ||\mathbf{m}_j - \mathbf{x}^{(n)}||^2]} \Rightarrow r^{(n)} = \text{softmax}(-\beta \left\{ ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2 \right\}_{k=1}^K)$$

$\beta \rightarrow \infty \Rightarrow$  Soft K-means becomes K-means. **Proof:** when beta is approaching infinity, the assignment step will be much more sensitive to the distance between each point and centers. Thus, if any center is slightly closer to the data point, it will immediately become the choice.

**EM Algorithm:**

Initialize the means  $\hat{\mu}_k$  and mixing coefficients  $\hat{\pi}_k$

Iterate until convergence:

E-step: Evaluate the responsibilities  $r_k^{(n)}$  given current parameters.

$$r_k^{(n)} = p(z^{(n)} = k | \mathbf{x}^{(n)}) = \frac{\hat{\pi}_k \mathcal{N}(\mathbf{x}^{(n)} | \hat{\mu}_k, I)}{\sum_{j=1}^K \hat{\pi}_j \mathcal{N}(\mathbf{x}^{(n)} | \hat{\mu}_j, I)} = \frac{\hat{\pi}_k \exp\{-\frac{1}{2} ||\mathbf{x}^{(n)} - \hat{\mu}_k||^2\}}{\sum_{j=1}^K \hat{\pi}_j \exp\{-\frac{1}{2} ||\mathbf{x}^{(n)} - \hat{\mu}_j||^2\}}$$

M-step: Re-estimate the parameters given the current responsibilities.

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N r_k^{(n)} \mathbf{x}^{(n)}, \hat{\pi}_k = \frac{N_k}{N} \text{ with } N_k = \sum_{n=1}^N r_k^{(n)}$$

Evaluate the log likelihood and check for convergence.

$\log p(D) = \sum_{n=1}^N \log(\sum_{k=1}^K \hat{\pi}_k \mathcal{N}(\mathbf{x}^{(n)} | \hat{\mu}_k, I)) \rightarrow$  When this converges, stop!

**What happened in the EM algorithm:** **1.** The maximum likelihood objective  $\sum_{n=1}^N \log p(\mathbf{x}^{(n)})$  was hard to optimize. **2.** However, the likelihood for  $D_{complete}$  is easy to compute.  $\log(D_{complete}) = \sum_{n=1}^N \sum_{k=1}^K I[z^{(n)} = k] (\log \mathcal{N}(\mathbf{x}^{(n)} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) + \log \pi_k)$ . **3.** We don't know about  $z^{(n)}$ , so we replaced  $I[z^{(n)} = k]$  with resiliities  $r_k^{(n)} = p(z^{(n)} = k | \mathbf{x}^{(n)})$ . That is, we replaced  $I[z^{(n)} = k]$  with its **expectation** under  $p(z^{(n)} | \mathbf{x}^{(n)})$ . Then, this is now easy to optimize, we can get  $\hat{\boldsymbol{\mu}}_k = \frac{1}{N_k} \sum_{n=1}^N r_k^{(n)} \mathbf{x}^{(n)}$ ,  $\hat{\pi}_k = \frac{N_k}{N}$  with  $N_k = \sum_{n=1}^N r_k^{(n)}$  (E-step)

**Both k-means and EM suffer from getting stuck at local minima.**

**Claim:** For EM algorithm with shared covariance  $\frac{1}{\beta} \mathbf{I}$ , it becomes soft k-means

algorithm. **Proof:** \_\_\_\_\_

**K-means vs EM:**

K-means	EM
Assignment Step: Assign each data point to the closest cluster. Refitting Step: Move each luster center to the average of the data assigned to it.	E-step: Compute the posterior probability over z given our current model. ( $r_k^{(n)} = p(z^{(n)} = k   \mathbf{x}^{(n)})$ ) M-step: Maximize the probability that it would generate the data it is currently responsible for.

## Matrix Factorization:

Can view the following as Matrix Factorization.

**PCA:**  $\hat{\mathbf{X}} = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ . SVD of the data matrix  $\mathbf{X} = \mathbf{Q} \mathbf{S} \mathbf{U}^T = \mathbf{Z} \mathbf{U}^T$ , with  $\mathbf{Z} = \mathbf{Q} \mathbf{S}$ .

The eigen-decomposition of  $\hat{\mathbf{X}}$  follows directly from the eigen-decomposition of  $\mathbf{X}^T \mathbf{X}$ :  $\hat{\mathbf{X}} = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \mathbf{U} \mathbf{S} \mathbf{Q}^T \mathbf{Q} \mathbf{S} \mathbf{U}^T = \mathbf{U} \left[ \frac{\mathbf{S}^2}{N} \right] \mathbf{U}^T$ . The SVD gives  $\mathbf{U}$ , which is equivalent to the learned basis of PCA. First K principal components corresponds first K columns of  $\mathbf{U}$ . Ultimately, PCA with K principal components finds the *optimal* rank-K approximation of  $\mathbf{X} \in \mathbb{R}^{N \times D}$ , in terms of error  $\|\mathbf{X}^T - \mathbf{U} \mathbf{Z}^T\|_F^2$ .

$\min \|\mathbf{X}^T - \mathbf{U} \mathbf{Z}^T\|_F^2$  over  $\mathbf{Z} \in \mathbb{R}^{N \times K}$ ,  $\mathbf{U} \in \mathbb{R}^{D \times K}$ . Note that the case K = D corresponds to the entire SVD of  $\mathbf{X}$ .

**Recommender Systems:**

Sparse ratings matrix:  $\mathbf{R}$ . Let the representation of user  $n$  be  $\mathbf{u}_n$  the representation of movie  $m$  be  $\mathbf{z}_m$ .  $n = 1, \dots, K$ . We want to enforce  $\mathbf{R} \approx \mathbf{U} \mathbf{Z}^T$ . It is impossible to do  $\min_{\mathbf{U}, \mathbf{Z}} \sum_{i,j} (R_{ij} - \mathbf{u}_i^T \mathbf{z}_j)^2$ , since most data entries  $R_{ij}$  are **missing!**

**Solution to missing data:** Let  $O = \{(n, m) : \text{every } (n, m) \text{ of matrix } \mathbf{R} \text{ is observed}\}$  Then we use square error loss to solve:  $\min_{\mathbf{U}, \mathbf{Z}} \sum_{(n,m) \in O} (R_{nm} - \mathbf{u}_n^T \mathbf{z}_m)^2$

Problem: this equation is very hard to minimize. It is non-convex. Solution is to use alternating least squares, A.k.a. ALS: fix  $\mathbf{Z}$  and optimize  $\mathbf{U}$ , then followed by fix  $\mathbf{U}$  and optimize  $\mathbf{Z}$ , and so on until convergence.

ALS for Matrix Completion algorithm

1. Initialize  $\mathbf{U}$  and  $\mathbf{Z}$  randomly
2. repeat until convergence
3.   for  $n = 1, \dots, N$  do
4.      $\mathbf{u}_n = \left( \sum_{m:(n,m) \in O} \mathbf{z}_m \mathbf{z}_m^T \right)^{-1} \sum_{m:(n,m) \in O} R_{nm} \mathbf{z}_m$
5.   for  $m = 1, \dots, M$  do
6.      $\mathbf{z}_m = \left( \sum_{n:(n,m) \in O} \mathbf{u}_n \mathbf{u}_n^T \right)^{-1} \sum_{n:(n,m) \in O} R_{nm} \mathbf{u}_n$

**K-means:** Can view K-means as a matrix factorization too. **1.** Stack 1-of-K vectors  $\mathbf{r}_i$  for assignments into a  $N \times K$  matrix  $\mathbf{R}$ , and stack the cluster centers  $\mathbf{m}_k$  into a matrix  $K \times D$  matrix  $\mathbf{M}$ . We want to let  $\mathbf{X} \approx \mathbf{R} \mathbf{M}$ , through reconstruction!

K-means distortion function in matrix form:  $\sum_{n=1}^N \sum_{k=1}^K r_k^{(n)} \|\mathbf{m}_k - \mathbf{x}^{(n)}\|^2 = \|\mathbf{X} - \mathbf{R} \mathbf{M}\|_F^2$ . We can do co-clustering so that  $\mathbf{X}$  will give a block structure.

**Sparse Coding:**  $\mathbf{x} \approx \sum_{k=1}^K s_k \mathbf{a}_k = \mathbf{A} \mathbf{s}$ , where  $\mathbf{s}$  is the sparse vector. **1.** We would like to choose  $\mathbf{s}$  to accurately reconstruct the image,  $\mathbf{x} \approx \mathbf{A} \mathbf{s}$  but encourage sparsity in  $\mathbf{s}$ .

**2.** We want to minimize the following function:  $\min_{\mathbf{s}} \|\mathbf{x} - \mathbf{A} \mathbf{s}\|^2 - \beta \|\mathbf{s}\|_1$

Here,  $\beta$  is a hyperparameter that trades off reconstruction error vs sparsity. We can learn a dictionary by optimizing both  $\mathbf{A}$  (aka the dictionary) and  $\{\mathbf{s}_i\}_{i=1}^N$  to trade off reconstruction error and sparsity

$$\min_{\{\mathbf{s}_i\}, \mathbf{A}} \sum_{i=1}^N \left( \|\mathbf{x}^{(i)} - \mathbf{A} \mathbf{s}_i\|^2 + \beta \|\mathbf{s}_i\|_1 \right), \text{ subject to } \|\mathbf{a}_k\|^2 \leq 1 \text{ for all } k$$

Can fit using ALS over  $\mathbf{A}$  and  $\mathbf{S}$ , just like anything mentioned before.

## Reinforcement Learning:

Goal: Find a policy  $\pi$  that maximizes the value function.

Optimal value function:  $Q^*(s, a) = \sup_{\pi} Q^{\pi}(s, a)$

Given  $Q^*$ , the optimal policy can be obtained as  $\pi^*(s) = \arg \max_a Q^*(s, a)$

The goal of an RL agent is to find a policy  $\pi$  that is close to optimal, i.e.,  $Q^{\pi} \approx Q^*$ .

The value function satisfy the following recursive relationship:

$$Q^{\pi}(s, a) = r(s, a) + \gamma \int_{\mathcal{S}} Q^{\pi}(s', \pi(s')) P(s' | s, a) ds' := (T^{\pi} Q^{\pi})(s, a)$$

$$(T^{\theta} Q)(s, a) := r(s, a) + \gamma \int_{\mathcal{S}} Q(s', \theta(s')) P(s' | s, a) ds'$$

$T^{\pi}$  is called the **Bellman operator**. We define the bellman optimality operator  $T^*$ :

$$(T^* Q)(s, a) := r(s, a) + \gamma \int_{\mathcal{S}} \max_{a' \in \mathcal{A}} Q(s', a') P(s' | s, a) ds'$$

Claim:  $Q^{\pi} = T^{\pi} Q^{\pi}$ ,  $Q^* = T^* Q^*$ . The solutions of these fixed-pt equations are unique.

**Value-based approaches** try to find a  $\hat{Q}$  such that  $\hat{Q} \approx T^* \hat{Q}$  (thus  $\hat{Q} \approx Q^*$ )

Then, we apply greedy policy of  $\hat{Q}$  in practice.

**Greedy policy of  $\hat{Q}$ :**  $\pi(s; \hat{Q}) := \arg \max_{a \in \mathcal{A}} \hat{Q}(s, a)$

**Finding the optimal value function using Value Iteration:**

Assume we know the model  $P$  and  $R$ , and the state-action space  $\mathcal{S} \times \mathcal{A}$  is small, Then we can use the bellman optimality equation. We start from an initial function  $Q_1$ .

For each  $k = 1, 2, \dots$ , apply  $Q_{k+1} \leftarrow T^* Q_k$ .

$$Q_{k+1}(s, a) \leftarrow r(s, a) + \gamma \int_{\mathcal{S}} \max_{a' \in \mathcal{A}} Q_k(s', a') P(s' | s, a) ds' \quad (\text{continuous})$$

$$Q_{k+1}(s, a) \leftarrow r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \max_{a' \in \mathcal{A}} Q_k(s', a') P(s' | s, a) \quad (\text{discrete})$$

**Claim:** The VI converges to the optimal value function. This is because of the contraction property of the Bellman (optimality) operator.

i.e.,  $\|T^* Q_1 - T^* Q_2\|_{\infty} \leq \gamma \|Q_1 - Q_2\|_{\infty}$ .

**Problem for VI: 1.** when state-action space is large, you cannot integrate exactly.

**2.** We often do not know the dynamics  $P$  and the reward function  $R$ , so we cannot calculate the Bellman operators. Solution  $\rightarrow$  Batch RL!

**Batch RL:** Suppose we are **only** given the following dataset.

1. A batch, i.d.  $D_N = \{(S_i, A_i, R_i, S'_i)\}_{i=1}^N$
2.  $(S_i, A_i) \sim v$  ( $v$  is a distribution over  $\mathcal{S} \times \mathcal{A}$ )
3.  $S'_i \sim P(\cdot | S_i, A_i)$ ; 4.  $R_i \sim R(\cdot | S_i, A_i)$

We can estimate  $Q \approx Q^*$  using the data given.

**Claim:**  $t_i = R_i + \gamma \max_{a' \in \mathcal{A}} Q(S'_i, a')$  is a noisy version of  $(T^* Q)(S_i, A_i)$ .

**Proof:**  $E[t_i | S_i, A_i] = E \left[ R_i + \gamma \max_{a' \in \mathcal{A}} Q(S'_i, a') | S_i, A_i \right] = r(S_i, A_i) + \dots = (T^* Q)(S_i, A_i)$

**The Algorithm of Batch RL:** Given the dataset  $D_N$  and action value function estimate  $Q_k$ , we firstly construct the dataset  $\{\mathbf{x}^{(i)}, t^{(i)}\}_{i=1}^N$  with  $\mathbf{x}^{(i)} = (S_i, A_i)$  and  $t^{(i)} = R_i + \gamma \max_{a' \in \mathcal{A}} Q(S'_i, a')$ . By the proof above we can treat the problem of estimating  $Q_{k+1}$  as a regression problem with noisy data. We solve a regression problem. We minimize the squared error.

$$Q_{k+1} \leftarrow \arg \min_{Q \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N \left| Q(S_i, A_i) - \left( R_i + \gamma \max_{a' \in \mathcal{A}} Q_k(S'_i, a) \right) \right|^2 =: \tilde{Q}_k$$

The policy of the agent is selected to be the **greedy** policy w.r.t. the final estimate of the value function: At state  $s \in \mathcal{S}$ , the agent chooses  $\pi(s; Q_k) \leftarrow \arg \max_{a \in \mathcal{A}} Q_k(s, a)$

This method is called **Approximate Value Iteration (AVI)**.

**Online Learning:** Motivation: RL problems are often interactive: The agent interacts with the environment and updates its knowledge of the world and its policy, with the goal of achieving as much rewards as possible.

**Q-Learning with epsilon-greedy policy:**

- Parameters:
  - Learning rate:  $0 < \alpha < 1$ : learning rate
  - Exploration parameter:  $\varepsilon$
- Initialize  $Q(s, a)$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$
- The agent starts at state  $S_0$ .
- For time step  $t = 0, 1, \dots$ ,
  - Choose  $A_t$  according to the  $\varepsilon$ -greedy policy, i.e.,

$$A_t \leftarrow \begin{cases} \arg \max_{a \in \mathcal{A}} Q(S_t, a) & \text{with probability } 1 - \varepsilon \\ \text{Uniformly random action in } \mathcal{A} & \text{with probability } \varepsilon \end{cases}$$

- Take action  $A_t$  in the environment.
- The state of the agent changes from  $S_t$  to  $S_{t+1} \sim P(\cdot | S_t, A_t)$
- Observe  $S_{t+1}$  and  $R_t$
- Update the action-value function at state-action  $(S_t, A_t)$ :

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_t + \gamma \max_{a' \in \mathcal{A}} Q(S_{t+1}, a') - Q(S_t, A_t) \right]$$

**RL miscellaneous:**

**Q:** What property does Transition probability follow? **S:** Markov property: the future depends on the past only through the current state

**Q:** Why do we need the discount factor? (Why cannot we always use the cumulative rewards, i.e.  $\gamma = 1$ ?) **S:** The cumulative reward is problematic for continuing tasks because the final time step would be  $T = \infty$ , and the return could easily be infinite.

**Q: VI vs AVI:** **S:** VI: Finding the optimal policy / value function when  $P$  and  $R$  known (planning problem) AVI: We are given dataset of size  $N$ , we can approximately perform VI using these data by regression, minimizing MSE.

**Q:** What is the **difference between Batch RL and online RL**? **S:** In batch RL, we have a batch of data coming from the previous interaction of the agent with the environment. This allowed us to use tools from the supervised learning literature. But RL problems are often interactive: the agent continually interacts with the environment, updates its knowledge of the world and its policy, with the goal of achieving as much rewards as possible.

**Q:** What is the difference between **exploration** and **exploitation** and how are they related in the Q-Learning algorithm? **S:** The RL agent should collect as much information about the environment as possible (exploration), while benefiting from the knowledge that has been gathered so far in order to obtain a lot of rewards (exploitation). In Q-Learning, the -greedy policy ensures that most of the time (probability  $1 - \varepsilon$ ) the agent exploits its incomplete knowledge of the world by choosing the best action, but occasionally (probability  $\varepsilon$ ) it explores other actions.