

Theory

Bayes classifier

g is the *classifier*.

$$\begin{aligned} g : \mathcal{X} &\rightarrow \mathcal{Y} \\ \mathbb{R}^d &\rightarrow \{0, 1\} \end{aligned}$$

To model the learning problem, we use the pair (X, Y) described by (μ, η) where μ is the probability measure:

$$\mu(A) = \mathbb{P}(X \in A)$$

And η is the regression of Y on X :

$$\eta(X) = \mathbb{P}(Y = 1 | X = x) = \mathbb{E}[Y | X = x]$$

η is also called the *a posteriori probability*.

The Bayes classifier is:

$$\begin{cases} 1 & \text{if } \eta(x) > 1/2 \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

Or, if \mathcal{Y} is $\{-1, 1\}$, we write the classifier as such: $g(x) = 2\mathbb{1}\{\eta(x) > 1/2\} - 1$.

Theorem:

For any classifier $g: \mathbb{R}^d \rightarrow \{0, 1\}$,

$$\mathbb{P}(g^*(X) \neq Y) \leq \mathbb{P}(g(X) \neq Y)$$

In other words, the Bayes classifier is theoretically **the best classifier**.

Proof: express $\mathbb{P}(g(X) \neq Y) - \mathbb{P}(g^*(X) \neq Y)$ in terms of dummies (use complementaries) and show that it is superior to 0.

Gradient descent

ℓ the loss function to minimize:

$$\theta_{t+1} = \theta_t - \alpha \nabla \ell$$

Algorithm 1 Global gradient descent

Loss function $\hat{L}_n(\hat{f}_\omega(x)) = \sum_{i=1}^n \ell(\hat{f}_\omega(x), y)$
 $E = 1000$
 $\epsilon = \text{small value}$
 $\omega_0 = \text{intial value in } t_0$
while $E > \epsilon$ **do**
 $\omega_{t+1} = \omega_t - \epsilon \sum_{i=1}^n \nabla_\omega \ell(\hat{f}_\omega(x), y)$
 Compute $E = L_n(\omega_{t+1})$
end while

Algorithm 2 Stochastic gradient descent

Loss function $\hat{L}_n(\hat{f}_\omega(x)) = \sum_{i=1}^n \ell(\hat{f}_\omega(x), y)$
 $E = 1000$
 $\epsilon = \text{small value}$
 $\omega_0 = \text{intial value in } t_0$
while $E > \epsilon$ **do**
 for $i = 1, \dots, n$ **do**
 $\omega_{t+1} = \omega_t - \epsilon \nabla_\omega \ell(\hat{f}_\omega(x_i), y_i)$
 end for
 Compute $E = L_n(\omega_{t+1})$
end while

Algorithm 3 Stochastic and random gradient descent

Loss function $\hat{L}_n(\hat{f}_\omega(x)) = \sum_{i=1}^n \ell(\hat{f}_\omega(x), y)$
 $E = 1000$
 $\epsilon = \text{small value}$
 $\omega_0 = \text{intial value in } t_0$
while $E > \epsilon$ **do**
 for $i = 1, \dots, n$ **do**
 Random draw of $i \in \{1, \dots, n\}$
 $\omega_{t+1} = \omega_t - \epsilon \nabla_\omega \ell(\hat{f}_\omega(x_i), y_i)$
 end for
 Compute $E = L_n(\omega_{t+1})$
end while

The main advantage of stochastic gradient is that it avoid computing the gradient descent on all the observations (greedy). However, in doing so, the gradient descent is subject to noise and can take longer to reach the optimum.

Note: the gradient is the derivation w.r.t ω

Note (stochastic and random gradient descent): the draw can be done with or without replacement.

Proof of the gradient descent formula:

<p>Taylor series by definition</p> $C(\theta_{new}) = C(\theta_{old}) + \nabla C \cdot [\theta_{new} - \theta_{old}] + ..$ <p>Here $\Delta\theta = [\theta_{new} - \theta_{old}]$, $\nabla C = \frac{\partial C}{\partial \theta}$</p> $C(\theta_{new}) \approx C(\theta_{old}) + \frac{\partial C}{\partial \theta} \cdot \Delta\theta \quad (1)$ <p>If we set $\Delta\theta = -\eta \frac{\partial C}{\partial \theta}$ with η a small positive learning rate equation (1) becomes:</p> $C(\theta_{new}) \approx C(\theta_{old}) + \frac{\partial C}{\partial \theta} \cdot \left(-\eta \frac{\partial C}{\partial \theta}\right) = C(\theta_{old}) - \eta \left(\frac{\partial C}{\partial \theta}\right) \left(\frac{\partial C}{\partial \theta}\right)$ <p>$C(\theta_{new}) \leq C(\theta_{old})$, since $\eta \left(\frac{\partial C}{\partial \theta}\right) \left(\frac{\partial C}{\partial \theta}\right)$ is always positive</p> <p>Conclusion : if we set $\Delta\theta = -\eta \frac{\partial C}{\partial \theta}$ it will decrease C</p>

(Idemia courses "DeepLearningTPT2018S1S2.pdf")

Learning rate optimization

Note:

- one epoch = one forward pass and one backward pass of all the training examples.
- batch size = the number of training examples in one forward/backward pass. The higher the batch size, the more memory space needed.
- number of iterations = number of passes, each pass using [batch size] number of examples. To be clear, one pass = one forward pass + one backward pass.

Learning rate decay

Simple idea: reduce the learning rate progressively.

E.g. $1/t$ decay:

$$\alpha_t = \frac{1}{(t+1)}$$

Momentum

Momentum is a method that helps accelerate SGD in the relevant direction and reduce oscillations.

General idea:

$$0 \leq \gamma \leq 1$$

$$M_{t_0} = x_0$$

$$M_{t_1} = \gamma M_{t_0} + x_1$$

$$M_{t_2} = \gamma M_{t_1} + x_2$$

Let's develop M_{t_2} to have a better view of momentum effect:

$$M_{t_2} = \gamma(\gamma M_{t_0} + x_1) + x_2 = \gamma^2 x_0 + \gamma x_1 + x_2$$

We can see that more importance is given to the most recent value (x_2) and the least to the past values.

In practice, $\gamma = 0.9$ gives good results.

Advantage: less dependent on noise

Adagrad - Adaptive Gradient Algorithm (2011)

Divide the learning rate by "average" gradient:

$$\theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{\sum_{i=0}^t (\nabla f_i)^2}} \nabla f$$

RMSProp - Root Mean Squared Propagation

Same as AdaGrad, but with an exponentially decaying average of past squared gradients.

$$\theta_t = \theta_{t-1} - \frac{\alpha}{\sigma_{t-1}} \nabla f$$

Where:

$$\sigma_t = \sqrt{\alpha(\sigma_{t-1})^2 + (1 - \alpha)(\nabla f_t)^2}$$

Adam - Adaptive moment estimation

Adam = RMS + momentum => use of a exponential decaying average of past squared gradients and past gradients

Bias-Complexity trade-off

ERM = Empiric Risk Minimization algorithm

\mathcal{H} = hypothesis class = all the classifiers that are considered

We can decompose the error of an $ERM_{\mathcal{H}}$:

$$L_{\mathcal{D}}(h_s) = \epsilon_{app} + \epsilon_{est}$$

- *Approximation error*: $\epsilon_{app} = \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$. This is the error done by the best predictor among those considered.

- *Estimation error*: $\epsilon_{est} = L_{\mathcal{D}}(h) - \epsilon_{app}$. This is the error difference from a used predictor and the best one.

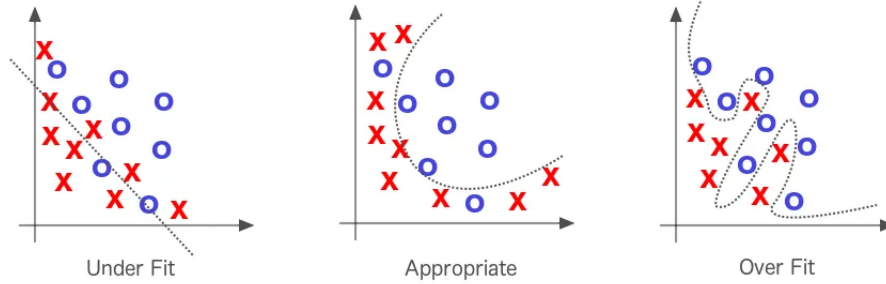
ϵ_{app} low => ϵ_{est} high => overfitting

ϵ_{est} low => underfitting

Overfitting

“Overfitting occurs when our hypothesis fits the training data "too well"”.

Understanding Machine Learning - From Theory To Algorithms.



Let S be a training set sampled according to the probability distribution \mathcal{D} . An algorithm A overfits if the difference between the *true risk* of its output $L_{\mathcal{D}}(A(S))$ and the *empirical risk* of its output $L_S(A(S))$ is large.

Note: Recall from the bias-complexity trade-off that overfitting is also when ϵ_{app} (error made by the best predictor) is low and thus ϵ_{est} (difference between the error of the best predictor and the error of the used predictor) is high.

Stability

An algorithm is stable if a small change in input gives a small change in output.

Let $S^{(i)}$ be a training set where we replace the i -th element by z' : $S^{(i)} = (z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m)$. We measure an effect of a small change comparing the losses $\ell(A(S^{(i)}), z_i)$ and $\ell(A(S), z_i)$ (z_i is the value to predict).

We note that the algorithm trained on $S^{(i)}$ doesn't observe z_i while the algorithm trained on S does. Intuitively, we should have $\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i) \geq 0$.

A is stable on average if $\mathbb{E}_{(S, z') \sim \mathcal{D}^{m+1}, i \sim U(m)} [\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)]$ is small.

Theorem:

$$\mathbb{E}_{S \sim \mathcal{D}} [L_{\mathcal{D}}(A(S)) - L_S(A(S))] = \mathbb{E}_{(S, z') \sim \mathcal{D}^{m+1}, i \sim U(m)} [\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)]$$

With $U(m)$ the uniform distribution over m .

Thanks to this theorem we have the following: A is a stable algorithm \iff second term is small \iff first term is small $\iff A$ does not overfit

Regularization

In this part we show that using regularization leads to a stable algorithm.

Let A be a regularized algorithm:

$$A(S) = \underset{\omega}{\operatorname{argmin}} (L_S(\omega) + \lambda \|\omega\|^2)$$

$$f_S : \omega \mapsto L_S(\omega) + \lambda \|\omega\|^2$$

Lemma:

- f_S is 2λ -strongly convex
- For any v : $f_S(v) - f_S(A(S)) \geq \lambda \|v - A(S)\|^2$ since $A(S)$ minimizes f_S

We have:

$$\begin{aligned} f_S(v) - f_S(u) &= L_S(v) + \lambda \|v\|^2 - (L_S(u) + \lambda \|u\|^2) \\ &= L_{S^{(i)}}(v) + \lambda \|v\|^2 - (L_{S^{(i)}}(u) + \lambda \|u\|^2) \\ &\quad + \frac{1}{m}(\ell(v, z_i) - \ell(v, z')) - \frac{1}{m}(\ell(u, z_i) - \ell(u, z')) \end{aligned}$$

We removed the loss computed on the additional observation z' and we added the loss computed on observation z_i .

$$\begin{aligned} f_S(v) - f_S(u) &= L_{S^{(i)}}(v) + \lambda \|v\|^2 - (L_{S^{(i)}}(u) + \lambda \|u\|^2) \\ &\quad + \frac{\ell(v, z_i) - \ell(u, z_i)}{m} + \frac{\ell(u, z') - \ell(v, z')}{m} \end{aligned}$$

Choosing $v = A(S^{(i)})$ and $u = A(S)$:

$$\begin{aligned} f_S(A(S^{(i)})) - f_S(A(S)) &= L_{S^{(i)}}(A(S^{(i)})) + \lambda \|A(S^{(i)})\|^2 - (L_{S^{(i)}}(A(S)) + \lambda \|A(S)\|^2) \\ &\quad + \frac{\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)}{m} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{m} \end{aligned}$$

If $A(S^{(i)})$ minimizes $L_{S^{(i)}}(\omega) + \lambda \|\omega\|^2$ (optimal coefficient):

$$L_{S^{(i)}}(A(S^{(i)})) + \lambda \|A(S^{(i)})\|^2 \leq L_{S^{(i)}}(A(S)) + \lambda \|A(S)\|^2$$

Thus,

$$f_S(A(S^{(i)})) - f_S(A(S)) \leq \frac{\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)}{m} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{m}$$

Thanks to the above lemma:

$$\lambda \|A(S^{(i)}) - A(S)\|^2 \leq \frac{\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)}{m} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{m} \quad (\text{I})$$

By definition, if $\ell(\cdot, z_i)$ is ρ -Lipschitz:

$$\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i) \leq \rho \|A(S^{(i)}) - A(S)\|$$

Similarly,

$$\ell(A(S^{(i)}), z') - \ell(A(S), z') \leq \rho \|A(S^{(i)}) - A(S)\|$$

Plugging these two equations in (I):

$$\lambda \|A(S^{(i)}) - A(S)\|^2 \leq \frac{2\rho \|A(S^{(i)}) - A(S)\|}{m}$$

$$\|A(S^{(i)}) - A(S)\| \leq \frac{2\rho}{\lambda m}$$

Using ρ -Lipschitz definition:

$$\frac{1}{\rho} (\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)) \leq \frac{2\rho}{\lambda m}$$

$$\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i) \leq \frac{2\rho^2}{\lambda m}$$

Since it is true for any S, z', i we conclude:

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) - L_S(A(S))] \leq \frac{2\rho^2}{\lambda m}$$

We have shown that a regularized algorithm is stable and thus does not overfit.

Supervised learning

Supervised learning aims at **finding a predictor when we have data with their labels**.

Perceptron

The Perceptron is the most basic algorithm for binary classification.

We want to estimate $f : \mathcal{X} \rightarrow \mathcal{Y}$ where $X \in \mathcal{X} \subset \mathbb{R}^p$ and $\mathcal{Y} = \{-1, 1\}$

Hyperplane

The classification is linear. We look for an hyperplane that separates the best the observations.

$$\mathcal{H} = \{x \in \mathbb{R}^p, \hat{f}_\omega(x) := \omega_0 + \sum_{i=1}^p \omega_i x_i = 0\}$$

We can also write:

$$\mathcal{H}_\omega : \omega^T x + \omega_0$$

Perceptron classifier: $x \mapsto \text{sign}(\hat{f}_\omega(x))$

Loss function

We want to find ω such that the loss function ℓ is minimised:

$\mathbb{E}[\ell(\hat{f}_\omega(x), y)]$ is the theoretic risk to minimize.

Recall that ERM consists of minimizing the empiric risk.

Empiric risk

The following is an example of loss function that can be used for the Perceptron problem:

$$\hat{L}_n(\hat{f}_\omega(x)) = -\sum Y_i(\omega_0 + \sum_{i=1}^p \omega_i x_i)$$

(when Y_i (target) and $\omega_0 + \sum_{i=1}^p \omega_i x_i$ (predicted) have different signs, the loss function increases)

Logistic regression

Logistic regression is used for binary classification.

It is quite similar to a simple linear regression in the sense that the objective is to find optimal weights ω to predict a variable. However, in the logistic regression we use a sigmoid function.

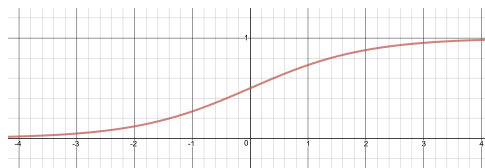
Rem: "logistic" because the logistic law has a sigmoid function as a repartition function.

Rationale behind the use of the sigmoid function:

We look for the *à posteriori* probability $\mathbb{P}(x|y=1) = \pi(x) = \hat{y}$.

The predicted variable \hat{y} is thus a probability.

The sigmoid function: $\sigma : z \rightarrow \frac{1}{1+e^{-z}}$ is well adapted because we want an output that is included in $[0, 1]$.



Classification function: $\hat{f}_\omega(x) = \sigma(\omega x)$ with a threshold

Loss function

If $y = 1$, we want $\sigma(\omega x)$ to be high $\Rightarrow 1 - \sigma(\omega x)$ should be low. The loss function should be increasing with $1 - \sigma(\omega x) = 1 - \frac{1}{1+e^{-\omega x}} = \frac{1}{1+e^{\omega x}}$. Equivalently, it should be increasing with $1 + e^{-\omega x}$.

More generally, the loss function is defined as such: $\ell(f_\omega, (x, y)) = \log(1 + e^{-y\omega x})$ (adding y in the expression allows to take into account cases when $y = 1$ and $y = -1$).

Recall that the log is a monotonic function.

Estimation

The advantage of the logistic loss function is that it is a convex function. Hence the ERM problem can be solved efficiently using standard methods.

Estimation is done using maximum likelihood. Maximum likelihood is finding the parameter that maximizes the probability to have a specific event (x_i, y_i) . We want to maximize the *à posteriori* probability that depends on x :

$$L(\omega, b) = \prod_{i=1}^n \pi(x_i)^{y_i} (1 - \pi(x_i))^{1-y_i}$$

This equation has no analytic solution. We use a numeric method to find the optimal parameters (see optimization algorithms).

See *Neural Network* section for more details on optimization.

Note: logistic regression is really a linear model since the objective is to find ω that is the slope of the line $\omega^T x + b$.

Linear discriminant analysis

We focus on the binary case, that is when $Y = +1$ or $Y = -1$.

These two conditional laws need to be gaussians with same covariance:

$X|Y = +1 \sim \mathcal{N}(\mu_+, \Sigma)$ with density f_+

$X|Y = -1 \sim \mathcal{N}(\mu_-, \Sigma)$ with density f_-

Let π_+, π_- be the simple probabilities $P(Y = +1), P(Y = -1)$

$$\mathbb{P}(Y = +1|X = x) = \frac{\mathbb{P}(Y=+1, X=x)}{\mathbb{P}(X=x)}$$

$$\mathbb{P}(Y = +1|X = x) = \frac{\mathbb{P}(X=x|Y=+1)\mathbb{P}(Y=+1)}{\mathbb{P}(X=x)}$$

$$\mathbb{P}(Y = +1|X = x) = \frac{f_+ \pi_+}{\mathbb{P}(X=x)}$$

$$\mathbb{P}(Y = +1|X = x) = \frac{f_+ \pi_+}{\mathbb{P}(X=x|Y=+1)\mathbb{P}(Y=+1) + \mathbb{P}(X=x|Y=-1)\mathbb{P}(Y=-1)}$$

$$\mathbb{P}(Y = +1|X = x) = \frac{f_+ \pi_+}{(f_+ \pi_+ + f_- \pi_-)}$$

Similarly,

$$\mathbb{P}(Y = -1|X = x) = \frac{f_- (1-\pi_+)}{\mathbb{P}(X=x)}$$

$$\mathbb{P}(Y = -1|X = x) = \frac{f_- (1-\pi_+)}{(f_+ \pi_+ + f_- \pi_-)}$$

The result shows us that we can express the two conditionnal probabilities in terms of conditionnal densities and "simple" probabilities (π_+, π_-).

Recall that multivariable gaussian density is: $f(x) = \frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$

In practice, μ_+, μ_-, π_+ and Σ are unknown. Thus we use empiric values:

$$\hat{\pi}_+ = m/n$$

$$\hat{\mu}_+ = \frac{1}{m} \sum \mathbb{1}_{\{y_i=+1\}} x_i$$

$$\hat{\mu}_- = \frac{1}{n-m} \sum \mathbb{1}_{\{y_i=-1\}} x_i$$

$$\hat{\Sigma} = \frac{1}{n-2} ((m-1)\hat{\Sigma}_+ + (n-m-1)\hat{\Sigma}_-)$$

$$\hat{\Sigma}_+ = \frac{1}{m-1} \sum \mathbb{1}_{\{y_i=+1\}} (x_i - \hat{\mu}_+)(x_i - \hat{\mu}_+)^T$$

$$\hat{\Sigma}_- = \frac{1}{n-m-1} \sum \mathbb{1}_{\{y_i=-1\}} (x_i - \hat{\mu}_-)(x_i - \hat{\mu}_-)^T$$

Classification

We predict class = 1 when $\mathbb{P}(Y = +1|X) > \mathbb{P}(Y = -1|X)$

$$\Rightarrow \frac{\mathbb{P}(Y=+1|X)}{\mathbb{P}(Y=-1|X)} > 1$$

$$\Rightarrow \log\left(\frac{\mathbb{P}(Y=+1|X)}{\mathbb{P}(Y=-1|X)}\right) > 0$$

Using previous conditional probability expressions, we end up with the following prediction rule:

$$\begin{cases} 1 & \text{if } x^T \hat{\Sigma}^{-1}(\hat{\mu}_+ - \hat{\mu}_-) > \frac{1}{2}\hat{\mu}_+^T \hat{\Sigma} \hat{\mu}_+ - \frac{1}{2}\hat{\mu}_-^T \hat{\Sigma} \hat{\mu}_- + \log(1 - m/n) - \log(m/n) \\ -1 & \text{otherwise} \end{cases} \quad (2)$$

$\hat{\mu}_+$, $\hat{\mu}_-$, $\hat{\pi}_+$ and $\hat{\Sigma}$ will be computed with *train data*.

x is the *test data*.

Note: $\hat{\Sigma}^{-1}(\hat{\mu}_+ - \hat{\mu}_-)$ is the **Fisher function** (Saporta).

Listing 1: LDA algorithm

```
class LDAClassifier():
```

```
    def fit(self, X, y):
```

```
        X_p = X[y == 1, :]
        X_m = X[y == -1, :]
```

```
        X_p_x1 = X_p[:, 0]
        X_p_x2 = X_p[:, 1]
        X_m_x1 = X_m[:, 0]
        X_m_x2 = X_m[:, 1]
```

```
        n = len(X)
```

```
        m = len(X_p)
```

```
        mean_p_x1 = np.mean(X_p_x1)
        mean_p_x2 = np.mean(X_p_x2)
        mean_p = np.array([mean_p_x1, mean_p_x2]) # mu_plus (estimated)
        cov_p = np.cov(np.transpose(X_p))
```

```
        mean_m_x1 = np.mean(X_m_x1)
        mean_m_x2 = np.mean(X_m_x2)
        mean_m = np.array([mean_m_x1, mean_m_x2]) # mu_minus (estimated)
        cov_m = np.cov(np.transpose(X_m))
```

```
        cov_est = (1/(n-2))*((m-1)* cov_p + (n-m-1)* cov_m) # sigma (estimated)
        inv_cov_est = np.linalg.inv(cov_est)
```

```
        a1 = np.dot(np.transpose(mean_p), inv_cov_est)
        a2 = np.dot(np.transpose(mean_m), inv_cov_est)
```

```
        # 2nd term in inequality
```

```
        self.alpha = 0.5*(np.dot(a1, mean_p) - 0.5*np.dot(a2, mean_m)) + np.log(1 - m/n) - np.log
```

```
        # 1st term in inequality
```

```
        self.beta = np.dot(inv_cov_est, mean_p - mean_m)
```

```
    return self
```

```

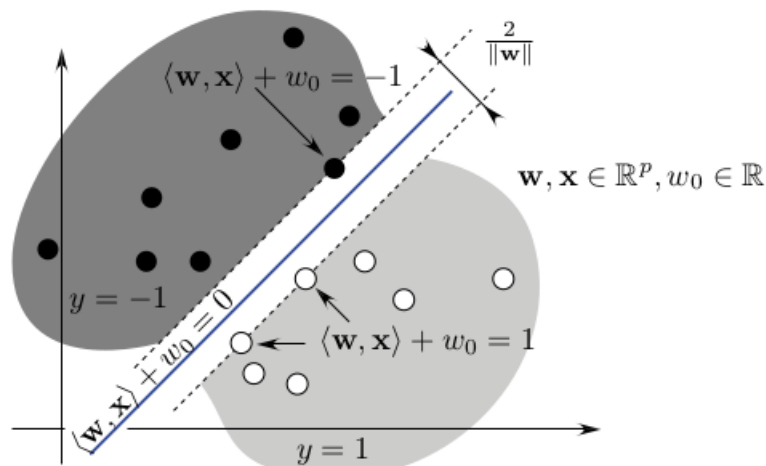
def predict(self, X):
    y_=[]
    for i in range(len(X)):
        X_pred = X[i]
        beta = np.dot(np.transpose(X_pred), self.beta)
        if (beta>self.alpha):
            Y_pred = 1
        else:
            Y_pred = -1
        y_.append(Y_pred)
    return np.array(y_)

```

SVM

Margin

This idea of SVM is to separate data as best as possible using a margin.



Three planes:

- $H_1 : \omega^T x + b = 1$
- $H : \omega^T x + b = 0$
- $H_{-1} : \omega^T x + b = -1$

Computing $H_1 - H_{-1}$ we get:

$$\begin{aligned}
 (x_1 - x_{-1})\omega^T &= 2 \\
 \Rightarrow \|x_1 - x_{-1}\| &= \frac{2}{\|\omega\|}
 \end{aligned}$$

The SVM problem starts with a **margin maximization**. We want to maximize the distance between x_1 and x_{-1} (to be double checked) and it is equivalent to minimizing $\|\omega\|$. Thus the optimization problem is written as such:

$$\begin{aligned}
 \min_{\omega, b} \quad & \frac{1}{2} \|\omega\|^2 \\
 \text{s.t.} \quad & y_i(\omega^T x_i + b) \geq 1 \quad i = 1, \dots, n
 \end{aligned}$$

Primal formulation

The above problem reflects the case when all data are perfectly separable, this is not true in practice. We thus add an error term ξ_i for each observation. This leads to the **primal formulation** of the problem:

$$\begin{aligned} \min_{\omega, b, \xi} \quad & \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\omega^T x_i + b) \geq 1 - \xi_i \quad i = 1, \dots, n \\ & \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

Note: the problem can be rewritten with the hinge loss function

$$\begin{aligned} y_i(\omega^T x_i + b) \geq 1 - \xi_i & \Rightarrow \xi_i \geq 1 - y_i(\omega^T x_i + b) \\ & \Rightarrow \xi_i \geq 1 - y_i(\omega^T x_i + b) \geq 0 \text{ since } \xi_i \geq 0 \\ & \Rightarrow \xi_i = \max(0, 1 - y_i(\omega^T x_i + b)) \\ & \Rightarrow \xi_i = (0, 1 - y_i(\omega^T x_i + b))_+ \\ & \Rightarrow \xi_i = \text{hinge}(f(x)) \text{ where } \text{hinge} \text{ is the Hinge loss function} \\ \text{hinge}(f(x)) &= (1 - yf(x))_+ \end{aligned}$$

$$\min_{\omega, b} \quad \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n (0, 1 - y_i(\omega^T x_i + b))_+$$

Lagrange

We can write this problem using Lagrange formulation, that is, integrating the constraints into the main formula:

$$\begin{aligned} \mathcal{L}(\omega, b, \xi, \alpha, \mu) &= \frac{1}{2} \omega^T \omega + C \sum \xi_i + \sum \alpha_i (1 - \xi_i - y_i(\omega^T x_i + b)) - \sum \mu_i \xi_i \\ \alpha_i, \mu_i &\geq 0 \end{aligned}$$

First order conditions:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \omega} = 0 & \Rightarrow \omega - \sum \alpha_i y_i x_i = 0 \Rightarrow \omega = \sum \alpha_i y_i x_i \\ \frac{\partial \mathcal{L}}{\partial b} = 0 & \Rightarrow -\sum \alpha_i y_i = 0 \\ \frac{\partial \mathcal{L}}{\partial \xi_i} = 0 & \Rightarrow C - \alpha_i - \mu_i = 0 \Rightarrow \alpha_i = C - \mu_i \\ \text{Since } \alpha_i, \mu_i &\geq 0 \text{ we have } C \geq \alpha_i \geq 0 \end{aligned}$$

Dual formulation

We can rewrite the problem using the first order conditions above:

$$\begin{aligned} \mathcal{L}(\omega, b, \xi, \alpha, \mu) &= \frac{1}{2} (\sum \alpha_i y_i x_i)^T (\sum \alpha_i y_i x_i) + C \sum \xi_i + \sum \alpha_i - \sum \alpha_i \xi_i - \sum \alpha_i y_i (\sum \alpha_i y_i x_i)^T x_i - \sum \alpha_i y_i b - \sum \mu_i \xi_i \\ &= -\frac{1}{2} (\sum \alpha_i y_i x_i)^T (\sum \alpha_i y_i x_i) + \sum (C - \alpha_i - \mu_i) \xi_i + \sum \alpha_i - \sum \alpha_i y_i b \\ &= -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum \alpha_i \end{aligned}$$

The problem is convex with lineary inequality constraints, we can apply the saddle point theorem.

Note: a point is saddle if it's a maximum w.r.t. one axis and a minimum w.r.t. another axis

The saddle theorem allows us to solve the problem $\min_{\omega} \max_{\alpha}$ as $\max_{\alpha} \min_{\omega}$

This leads to the problem in its **dual formulation**:

$$\begin{aligned}
\max_{\alpha} \quad & -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum \alpha_i \\
\text{s.t.} \quad & 0 \leq \alpha_i \leq C \quad i = 1, \dots, n \\
& \sum \alpha_i y_i = 0 \quad i = 1, \dots, n
\end{aligned}$$

Using the above expression of ω (optimal condition), the classification function is:

$$f(x) = \text{sign}(\sum \alpha_i y_i x_i^T x + b)$$

Kernel

Kernels are used when separation is non linear.

Recall the primal formulation:

$$\begin{aligned}
\min_{\omega, b, \xi} \quad & \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i \\
\text{s.t.} \quad & y_i (\omega^T x_i + b) \geq 1 - \xi_i \quad i = 1, \dots, n \\
& \xi_i \geq 0 \quad i = 1, \dots, n
\end{aligned}$$

When separation is non linear, we set ϕ as a non linear transformation. The constraint becomes:
 $y_i (\omega^T \phi(x_i) + b) \geq 1 - \xi_i \quad i = 1, \dots, n$

Dual formulation is now:

$$\begin{aligned}
\max_{\alpha} \quad & -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j) + \sum \alpha_i \\
\text{s.t.} \quad & 0 \leq \alpha_i \leq C \quad i = 1, \dots, n \\
& \sum \alpha_i y_i = 0 \quad i = 1, \dots, n
\end{aligned}$$

The classification becomes:

$$f(x) = \text{sign}(\sum \alpha_i y_i \phi(x_i)^T \phi(x) + b)$$

To classify a new point, we thus need to be able to compute $\phi(x_i)^T \phi(x)$.

Kernel trick: there is no need to know an explicit expression of ϕ (i.e. knowing the coordinates of points in new set) since we are only looking at distances and angles, that is scalar product.

Kernel functions implement those scalar products: $K(x, x') = \phi(x)^T \phi(x')$ where ϕ is a transformation function into a Hilbertian set $\phi: \mathcal{X} \rightarrow \mathcal{F}$

Note: an Hilbertian is a set with scalar product: $\mathcal{F} = (\mathcal{H}, \langle \cdot, \cdot \rangle)$

Most used kernels:

Linear kernel: $K(x, x') = x^T x'$ (we often call this setup a "no-kernel SVM")

Polynomial kernel: $K(x, x') = (x^T x' + c)^d$

Gaussian kernel: $K(x, x') = \exp(-\gamma \|x - x'\|^2)$

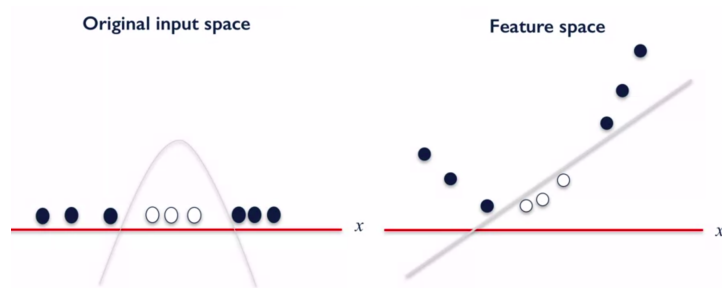
The optimisation problem can be written with kernel:

$$\begin{aligned}
\max_{\alpha} \quad & -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(x_i, x_j) + \sum \alpha_i \\
\text{s.t.} \quad & 0 \leq \alpha_i \leq C \quad i = 1, \dots, n \\
& \sum \alpha_i y_i = 0 \quad i = 1, \dots, n
\end{aligned}$$

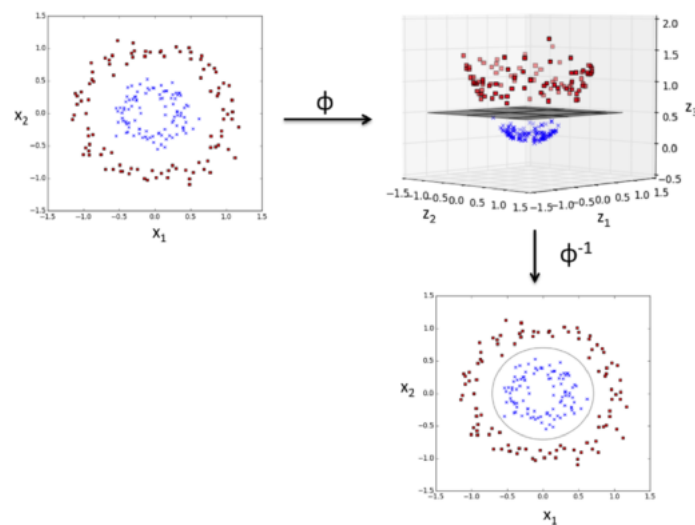
Summary

SVM allows to find complex non linear separations in transforming the problem into a higher dimension where data are linearly separable.

From 1D to 2D:



From 2D to 3D:



Neural network

Logistic regression with a neural network mindset

Based on the example from coursera ([deeplearning.ai](https://www.coursera.org/learn/deeplearning))

Let us take the example of an image that we want to classify in a **binary** way: man/woman

The picture is vectorized as a vector of pixels : $\begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$

We use a regression to predict if it's a man/woman:

$$y = \omega^T x + b$$

Note: x are all the pixels of **one** image.

We want a probability in output (if it's ≥ 0.5 then we say it's a man).

We thus want the output to be $\hat{y} = \sigma(\omega^T x + b) = \mathbb{P}(y|x) \in [0, 1]$

(see regression part to get more details on the sigmoid)

Now since it's a binary classification, we want the y (real value) to be 0 or 1.

Thus, the loss function is:

$$\mathcal{L}(y, \hat{y}) = -y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

The cost function is the empiric loss on all examples:

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^i)$$

Forward propagation

$$x_1, x_2, \omega_1, \omega_p, b \rightarrow z = \omega_1 x_1 + \omega_2 x_2 + b \rightarrow \hat{y} = a = \sigma(z) \rightarrow \mathcal{L}(a, y)$$

- First arrow: regression
- Second arrow: probability
- Third arrow: error

Backward propagation

The idea is: with the error computed on the last step, we go backward in order to correct the parameters ω and b .

$$x_1, x_2, \omega_1, \omega_p, b \leftarrow z = \omega_1 x_1 + \omega_2 x_2 + b \leftarrow \hat{y} = a = \sigma(z) \leftarrow \mathcal{L}(a, y)$$

Example: we want to find ω_1 that minimizes the cost function:

$$\frac{d\mathcal{L}}{d\omega_1} = "d\omega_1" = \frac{d\mathcal{L}}{da} \frac{da}{dz} \frac{dz}{d\omega_1} = \dots = (a - y)x_1 = dzx_1$$

Steps:

We compute all the derivatives, then we apply the gradient descent

Listing 2: Gradient descent (logistic regression with a NN mindset)

```
for i in range(num_iterations):  
  
    # Cost and gradient calculation  
    grads, cost = propagate(w, b, X_train, Y_train) # propagation on ALL the training sample  
  
    # Retrieve derivatives from grads  
    dw = grads["dw"]  
    db = grads["db"]  
  
    # update parameters  
    w = w - learning_rate * dw
```

```
b = b - learning_rate * db
```

```
# Record the costs
costs.append(cost)
```

Listing 3: Propagation (logistic regression with a NN mindset)

```
def propagate(w, b, X, Y):
```

```
    m = X.shape[1]
```

```
    # FORWARD PROPAGATION (FROM X TO COST)
```

```
    A = sigmoid(np.dot(w.T, X) + b)
```

```
    cost = (- 1 / m) * np.sum(Y * np.log(A) + (1 - Y) * (np.log(1 - A)))
```

```
    # BACKWARD PROPAGATION (TO FIND GRAD)
```

```
    dw = (1/m)*np.dot(X, (A-Y).T)
```

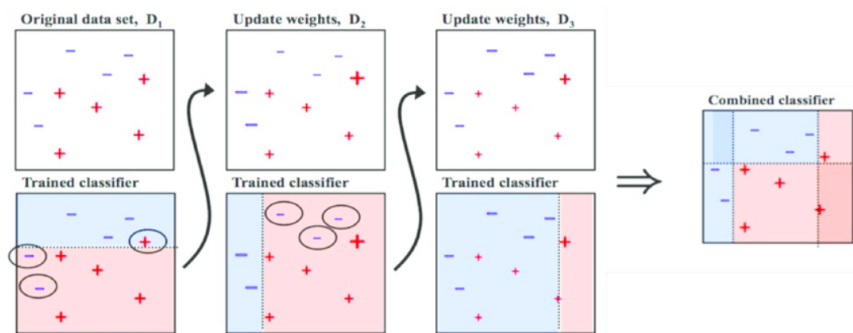
```
    db = (1/m)*np.sum(A-Y)
```

AdaBoost

Boosting is an algorithmic paradigm addressing two major issues in machine learning:

- It optimizes the bias-complexity trade-off. The learning starts with a basic class (large approximation error) and as it progresses the hypothesis class becomes more complex.
- It allows to find predictors that are usually computationally infeasible to find.

Main idea: weak learners are "boosted" to become stronger altogether.



Weak learner (or γ -weak-learner): it's an **algorithm** returning a function h such that $L_{\mathcal{D}}(h) \leq 1/2 - \gamma$. In other words, it returns a simple binary predictor that does slightly better than a random guess.

We note:

- Final predictor = weighted sum of weak predictors
- More weights are given to observations that gave wrong prediction. In doing so, the classifier of the next round will focus on these observations. Warning: to see this, focus on the variation of $D_i^{(t+1)}$ and not just w_t .

Theorem: the training error of the output hypothesis decreases **exponentially fast** with the number of boosting rounds.

Algorithm 4 AdaBoost

Input: training set $S = (x_1, y_1), \dots, (x_m, y_m)$, weak learner WL, number of rounds T

Initialize $D^{(1)} = (1/m, \dots, 1/m)$ // same weights for all observations

for $t = 1, \dots, T$ **do**

 invoke weak learner $h_t = WL(D^{(t)}, S)$

 compute $\epsilon_t = \sum_{i=1}^m D_i^{(t)} \mathbb{1}_{[y_i \neq h_t(x_i)]}$ // misclassification error

 let $\omega_t = \frac{1}{2} \log(\frac{1}{\epsilon_t} - 1)$ // $\epsilon_t < 1$

 update $D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\omega_t y_i h_t(x_i))}{\sum_{j=1}^m D_j^{(t)} \exp(-\omega_t y_j h_t(x_j))}$ for all $i = 1, \dots, m$

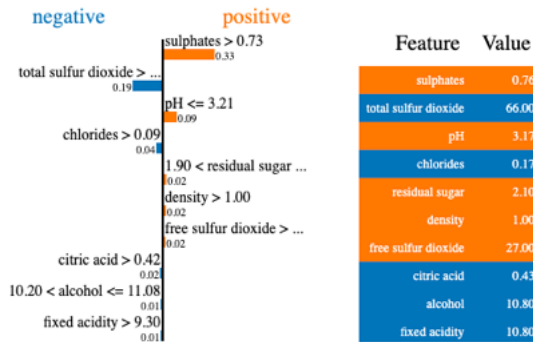
end for

Output: the hypothesis $h_s(x) = \text{sign}(\sum_{t=1}^T \omega_t h_t(x))$

LIME

Local Interpretable Model-agnostic Explanations

LIME is a model allowing to explain why an algorithm made a specific decision on a specific observation. It is model-agnostic since the method can explain decisions without understanding how the classifier works. By "explaining" we mean displaying feature importance for the decision. Typical output is:



On this picture, we see the features contributing to wine quality: sulphates positively contributes to the wine quality while sulfur dioxide contributes negatively.

Model

The main idea behind LIME is to perturb data and learn an interpretable model locally.

Input: observation we want to explain x_0 .

Step 1: local perturbation of x_0

A neighborhood is created around x_0 . By default, LIME uses a Gaussian sampling: it generated variables from a normal distribution. Then it de-standardizes generated samples using mean and standard deviation from the observation to explain.

Note: by default, LIME generates a neighborhood of 5000 samples.

Step 2: weight computation

Weights π_{x_0} are computed according to the distance between x_0 and its neighbors. LIME uses an exponential smoothing kernel:

$$K(x, y) = e^{-\frac{\|x-y\|_2^2}{\sigma^2}}$$

Reminder: $\|x - y\|_2 = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$

Step 3: local classification

A linear model is used to classify samples around x_0 . A linear model is used because it's an *interpretable* model (such as decision trees).

Loss minimization to explain x_0 :

$$\xi(x_0) = \underset{g \in G}{\operatorname{argmin}} \mathcal{L}(f, g, \pi_{x_0}) + \Omega(g)$$

Weights π_{x_0} are taken into account in the loss function.

$\Omega()$ is a complexity measure. By default, LIME uses a Ridge regression, thus $\Omega()$ is the 2-norm. It would also make sense to use a lasso regularized regression (1-norm) in order to allow zero coefficients (and thus have variable selection).

Step 4: display explanations

Feature importance is displayed looking at regression coefficients.

Unsupervised learning

Unsupervised learning aims at **learning some underlying hidden structure of the data when we don't have the labels**.

Unsupervised models can be used as a pre step for supervised learning, e.g.:

- reduce the training sample (dimensionality reduction: forward selection, PCA, autoencoders)
- give output for unlabeled data (clustering, autoencoders)
- grow the training sample (generative models)

Expectation-Maximization (EM) in the case of GMM (Gaussian Mixture Model)

(for more details, see document *gmm.pdf* in Cloud folder)

GMM problem aims at estimating parameters of a sample distribution. ("Generative models" p. 295 Understanding Machine Learning).

A GMM sample is composed of j Gaussian variables (*clusters*) distributed with proportions (π_1, \dots, π_k) ($\sum \pi_i = 1$)

We can write:

$$X \sim \mathcal{N}(\mu_Z, \Sigma_Z) \quad \text{with } Z \sim \pi$$

π is not really a law but more the proportions of each Gaussian categories.

Thus, X has a density which is a weighted-average of all Gaussian densities:

$$p_{\theta}(x) = \sum_{j=1}^k \pi_j f_j(x) \quad (*)$$

Estimation

We want to estimate $\theta = (\pi, \mu, \Sigma)$ where:

$$\pi = (\pi_1, \dots, \pi_k), \mu = (\mu_1, \dots, \mu_k), \Sigma = (\Sigma_1, \dots, \Sigma_k)$$

To do so, we use the maximum likelihood method (product of densities across all samples):

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i)$$

$$l(\theta) = \log(\prod_{i=1}^n p_{\theta}(x_i)) = \sum_{i=1}^n \log(p_{\theta}(x_i))$$

We thus need to find $\argmax(l(\theta))$

Problem: the likelihood function is not convex!

The expectation-maximization problem is used when we have *latent variables* (= variables for which we don't know their associated distribution).

Let $z = (z_1, \dots, z_k)$ be the vector of latent variables. We can express the density (*) as a joint function with respect to z :

$$p_{\theta}(x, z) = p_{\theta}(z) p_{\theta}(x|z)$$

$$l(\theta, z) = \dots = \sum (\log \pi_{z_i}) + \sum (\log f_{z_i}(x_i))$$

A classic optimization (in case of Gaussians) give us empirical values as solutions e.g. $\hat{\pi}_j = \frac{n_j}{n}$

Problem: we don't know j !

We will thus use the *expected* log-likelihood method.

Let us find another expression of the likelihood:

$$p_{\theta}(x, z) = p_{\theta}(x) p_{\theta}(z|x)$$

As seen previously: $p_{\theta}(x, z) = \prod \pi_{z_i} f_{z_i}(x_i)$

$$p_{\theta}(z|x) = \prod p_{\theta}(z_i|x_i) = \frac{\prod \pi_{z_i} f_{z_i}(x_i)}{p_{\theta}(x_i)} \propto \prod \pi_{z_i} f_{z_i}(x_i)$$

Given an initial parameter θ_0 , the *expected* log-likelihood is written as such:

$$\mathbb{E}_{\theta_0}[l(\theta; z)] = \sum p_{\theta_0}(z|x) l(\theta; z)$$

$$\mathbb{E}_{\theta_0}[l(\theta; z)] = \sum_j \sum_i p_{ij} (\log \pi_j + \log f_j(x_i))$$

We now have an expression that doesn't depend on z but only on p_{ij} and we know that $n_j = \sum_i p_{ij}$

K-means

(see kmeans.pdf from OneDrive folders for more details)

Objective: group data into k clusters so that samples in the same cluster are close to each other w.r.t. the Euclidean distance.

The cost function minimization is written as such:

$$\underset{C_1, \dots, C_k; \mu_1, \dots, \mu_k}{\operatorname{argmin}} \quad \sum_{j=1}^k \sum_{i \in C_j} \|x_i - \mu_j\|^2$$

Where μ_j is the mean, also called gravity center or cluster center:

$$\mu_j = \frac{1}{|C_j|} \sum_{i \in C_j} x_i$$

Algorithm 5 K-means

```

Input:  $x_1, \dots, x_n$ 
Output: clusters  $C_1, \dots, C_k$ 
Random values for  $\mu_1, \dots, \mu_k$ 
while no convergence do
    // Step 1: Update clusters
     $C_1, \dots, C_k \leftarrow \emptyset$ 
    for  $i = 1$  to  $n$  do
         $j \leftarrow \operatorname{argmin}_l \|x_i - \mu_l\|$ 
         $C_j \leftarrow C_j + \{i\}$  // We add observation  $i$  to the cluster  $C_j$ 
    end for
    // Step 2: Update cluster centers
    for  $j = 1$  to  $k$  do
         $\mu_j \leftarrow 0$ 
         $n_j \leftarrow 0$ 
        for  $i$  in  $C_j$  do // We loop on all observations of each cluster
             $\mu_j \leftarrow \mu_j + x_i$ 
             $n_j \leftarrow n_j + 1$ 
        end for
         $\mu_j \leftarrow \mu_j / n_j$ 
    end for
end while
```

The quality of the clustering strongly depends on the initial center values. This is why the algorithm is generally run multiple times for different initial values. The best clustering (i.e., that of minimum cost) is returned.

K-means++

To improve the quality of the clustering, we choose the initial cluster centers far from each other:

- select the first cluster center uniformly at random among the n data samples
- select the following cluster centers at random, with a probability proportional to the square distance to the closest current cluster center

Listing 4: K-means++ initial centers selection

```

centers = []
centers.append(X[np.random.randint(X.shape[0])]) # initial center = one random sample
distance = np.full(X.shape[0], np.inf) # a vector (n,1) with only infinity terms
for j in range(1, self.n_clusters):
    distance = np.minimum(np.linalg.norm(X - centers[-1], axis=1), distance) # size (n,1);
# distance = the smallest distance associated with
# the last added center
    p = np.square(distance) / np.sum(np.square(distance)) # probability vector [p1,...,pn]
# the highest probability in p is associated
# with the biggest distance w.r.t the last added center
    sample = np.random.choice(X.shape[0], p = p) # one sample is
# selected according to probabilities

centers.append(X[sample])

```

Note: this problem is called *NP-hard problem*. It means that its complexity is at least equal to the complexity of an NP-problem

NP-problem: a problem is NP if it can be determined by a non-deterministic Turing machine in polynomial time. Intuitively, a problem is NP if we can quickly verify if one candidate is a solution of the problem. E.g. "travelling salesman problem" = let d be a distance and n be a number of cities. Is there an itinerary with distance $\geq d$ stopping by every city? -> easy to check...

Turing machine (1936)



"non-deterministic turing machine": itinerary can be represented by a tree...

Local Outlier Factor

Local Outlier Factor is an unsupervised method used in anomaly detection. It consists of comparing local density of train observations VS local density of test observations.

Reachability distance

$reachability-distance_k(A, B) = \max\{k-distance(B), d(A, B)\} = \text{reachability of A from B.}$

$k-distance(B)$: distance from B to its kth nearest neighbor.

The reachability distance of A from B is *at least* the distance between A and B or *at least* the distance of B's neighbor.

When A is very far from B, it's simply the distance between the two points.
When A is very close to B, it's the distance between B and its neighbor.

The distance can be computed using different metrics: Euclidean distance, Mahalanobis distance, etc.

Local reachability density

$$lrd_k(A) = \frac{1}{\sum_{B \in N_k(A)} reachability-distance_k(A, B) / |N_k(A)|}$$

It's the inverse of the average of reachability-distances of A from B.

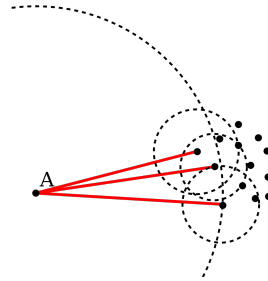
When A is very far from its neighbors: sum of reachability distances is high => local reachability density is small.

Local Outlier Factor

LOF computation consists of comparing the local densities of a point VS its neighbors.

$$LOF_k(A) := \frac{\frac{\sum_{B \in N_k(A)} lrd_k(B)}{lrd_k(A)}}{|N_k(A)|}$$

$LOF_k(A) > 1$: A is an outlier. Local reachability density of A is small compared to its neighbors.
 $LOF_k(A) < 1$: A is an inlier.



On this figure, $k = 3$. We can see that the reachability distance of A from its neighbors is high (red segments). The local reachability density of A will thus be **low**.

On the contrary, the local reachability densities of its neighbors is **high** because each neighbor can be easily reached from their own neighbors.

As a result, LOF would be high so A is an outlier.

sklearn algorithm

To score an observation, *fit* simply memorizes the train observations (same as in knn).

score_samples first finds the k-nearest neighbors from the train set thanks to the given distance metric. It then computes the local outlier factor for each test observation comparing the test observation local density with its closest k-neighbors local densities in the train set.

Variational Auto-Encoder

Variational autoencoders are a combination of three things:

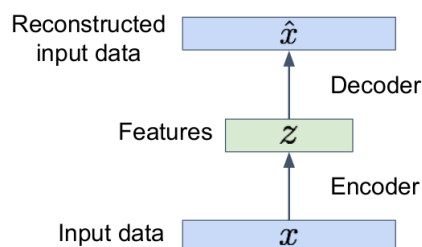
1. Autoencoders
2. Variational Approximation & Variational Lower Bound
3. "Reparameterization" Trick

1. Autoencoders

Autoencoders are used to extract features from unlabeled training data. They are new methods for **dimensionality reduction** and part of neural networks branch.

Note: autoencoders can be used to replace older dimensionality reduction methods such as PCA for several reasons:

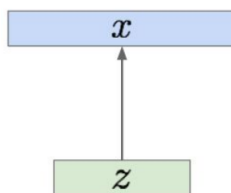
- For very large data sets that can't be stored in memory, PCA will not be able to be performed. The autoencoder construction using keras can easily be batched resolving memory limitations.
- PCA is restricted to linear separation while autoencoders are capable of modelling complex non linear functions.



Learning can be done using a loss function such as $\|x - \hat{x}\|^2$. In a similar way than neural network, optimization is typically done with backpropagation.

2. Variational Approximation & Variational Lower Bound

We assume x is generated from unobserved (latent) z :



Practical example: x can be seen as images and z as the main attributes (orientation, colors, etc.)

$x \sim p_{\theta^*}(x|z)$ where $p_{\theta^*}(x|z)$ is called *true conditional*
 $z \sim p_{\theta^*}(z)$ where $p_{\theta^*}(z)$ is called *true prior*

Objective: estimating $p_{\theta}(x)$. We thus need to estimate θ^* .

We can do it through maximum likelihood. The marginal density is $p_\theta(x) = \int p_\theta(x|z)p_\theta(z)dz$

Note: a marginal likelihood function is a likelihood function in which some parameter variables have been marginalized. Marginalization consists in summing over the possible values of one variable in order to determine the contribution of another. E.g., $\mathbb{P}(X) = \sum_y \mathbb{P}(X, Y = y)$ or in continuous probabilities $p(x) = \int p(x, y)dy$. Also, if we don't know the joint probability, we can express this using conditional probabilities: $p(x) = \int p(x|y)p(y)dy$

Problem: impossible to compute $p(x|z)$ for every z (**computationally too expensive**) => problem is said **intractable**

Solution: use another encoder learning $q_\phi(z|x)$ that approximates $p_\theta(z|x)$

$$\begin{aligned}
\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] && (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\
&= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) p_\theta(z)}{p_\theta(z | x^{(i)})} \right] && (\text{Bayes' Rule}) \\
&= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] && (\text{Multiply by constant}) \\
&= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] && (\text{Logarithms}) \\
&= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))
\end{aligned}$$

$\mathbf{E}_z [\log p_\theta(x^{(i)} | z)]$: we can estimate this term through sampling

$D_{KL}(q_\phi(z|x^{(i)}) || p_\theta(z))$: differentiable term

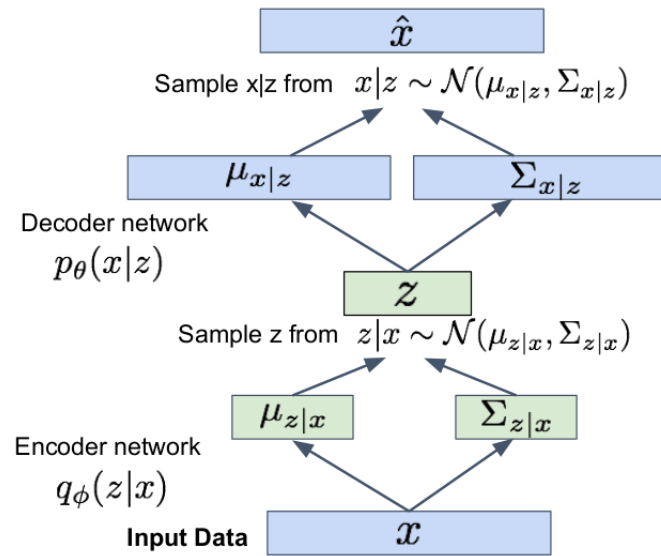
$D_{KL}(q_\phi(z|x^{(i)}) || p_\theta(z|x^{(i)}))$: $p(z|x)$ intractable but we know that $D_{KL} \geq 0$

Let $\mathcal{L}(x^{(i)}, \theta, \phi) = \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z|x^{(i)}) || p_\theta(z)) =$ **tractable lower bound** that we can optimize

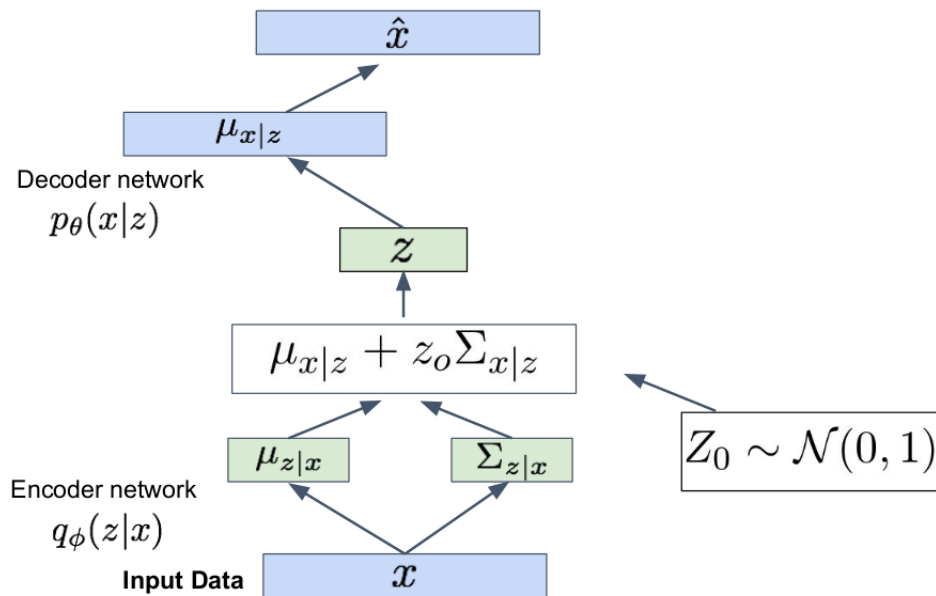
We know that $p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$ since $D_{KL}(q_\phi(z|x^{(i)}) || p_\theta(z|x^{(i)})) \geq 0$

Thus the maximum likelihood problem becomes: $\theta^*, \phi^* = \argmax_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$

We can minimize $D_{KL}(q_\phi(z|x^{(i)}) || p_\theta(z))$ making posterior distribution close to prior. To do so, we make encoder network predicting $\mu_{z|x}$ and $\Sigma_{z|x}$ and then we sample $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$



Problem: sampling $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$ and $x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$ is not differentiable (why?).
 \Rightarrow we use **reparametrization trick**: we sample $z_0 \sim \mathcal{N}(0, 1)$ to have $z = \mu_{x|z} + z_0 \Sigma_{x|z} \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$



Optimization through forward and backward propagation! **Reinforcement learn-**
ing

Reinforcement learning is inspired on human logic: we learn which strategy to take thanks to rewards we receive.

Reinforcement learning uses Markov Decision Processes (MDP).

A Markov Decision Process is defined by:

- an initial state s_0
- the reward distribution $r_t \sim p(r|s_t, a_t)$ (stochastic)
- the transition probabilities, $s_{t+1} \sim p(s|s_t, a_t)$ (stochastic)

=> **MDP: we know the state and the reward from the previous state only.**

A *policy* is an action for each state, for a given MDP.

=> **a policy can be seen as a strategy: we know where to go at each state.** In practice the policy is actually the **transition probability matrix**.

The *value function* is the gain we earn at a state, for a specific policy: $\forall s, V_\pi(s) = \mathbb{E}_\pi[G|s_0 = s]$

=> **The expected gain takes into account the transition probability.**

The gain is calculated as such: $G = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots = \sum_t \gamma^t r_t$ where γ is the *discount factor*. V can also be written $V_\pi(s) = \mathbb{E}_\pi[r_0 + \gamma V(s_1)|s_0 = s]$. This expression is called **Bellman equation**.

We can find the best policy:

$$\forall s, \pi^*(s) = a^* \in \operatorname{argmax}_a \mathbb{E}_\pi[r_0 + \gamma V_*(s_1)|s_0 = s, a_0 = a]$$

Where V_* is solution of the **Bellman optimality equation**:

$$\forall s, V(s) = \max_a \mathbb{E}[r_0 + \gamma V(s_1)|s_0 = s]$$

Online estimation

The expected value function is approached using empiric estimator (sum).

Two ways to do it:

- Monte-Carlo update: if we can memorize all the paths, we update the sum at each step $S \leftarrow x_t$.

At the end we compute the mean $X \leftarrow \frac{S}{t}$

- TD-learning: we update the value function using temporal differences $\forall s, V(s_t) \xleftarrow{\alpha} r_t + \gamma V(s_{t+1})$

Where $X \xleftarrow{\alpha} x_t \Leftrightarrow X = X + \alpha(x_t - X)$ (α is usually $1/t$)

Online control

Recall that value function is $\forall s, V_\pi(s) = \mathbb{E}[G|s_0 = s]$

Unknown model \Leftrightarrow "unknown strategy" \Leftrightarrow we don't know the rewards in advance (we need to learn *online*).

=> we cannot compute the expectation \mathbb{E}

=> instead of the value function, we will estimate the **value-action function** thanks to the Bellman equation:

$$\forall s, a, \quad Q_\pi(s, a) = \mathbb{E}_\pi[r_0 + \gamma Q(s_1, a_1)|s_0 = s, a_0 = a]$$

=> for an initial state and a fixed action, we can compute the Q function

=> the value-action function can be seen as the expected (= > use of transition matrix (probabilities)) gain we get at a state when doing a specific action

Note: the policy π takes account only from next state s_1

We want to *control* the policy and find the optimal one: $\pi^*(a) = a^* \in \operatorname{argmax}_a Q_*(s, a)$

The optimal Bellman equation becomes:

$$\forall s, a, \quad Q_\pi(s, a) = \mathbb{E}_\pi[r_0 + \gamma \max_{a'} Q(s_1, a') | s_0 = s, a_0 = a]$$

How to choose initial state a_0 ?

-> Pure exploitation: we find the best action knowing the current system $\pi(s) \leftarrow \operatorname{argmax}_a Q(s, a)$

=> problem: Q is estimated, thus we have a chance to miss good actions

-> Pure exploration: $\pi(s) \leftarrow \text{random}$

=> problem: we can waste time on bad actions and thus have bad estimation quality

Solution: trade-off exploitation/exploitation (SARSA, Q-learning, ...).

SARSA

This algorithm is based on ϵ -greedy algorithm.

$$\pi(s) \leftarrow \begin{cases} a^* \in \operatorname{argmax}_a Q(s, a) & \text{with probability } 1 - \epsilon \\ \text{random} & \text{with probability } \epsilon \end{cases} \quad (3)$$

=> with a small probability (ϵ is usually lower than 5%), we do random actions. Otherwise we take the best action of the known ones.

Estimation is done through *TD-learning* (temporal differences - see Online Estimation):

$$\forall t, Q(s_t, a_t) \stackrel{\alpha}{\leftarrow} r_t + \gamma Q(s_{t+1}, a_{t+1})$$

Listing 5: SARSA

```
# action is identified with new state (after move) except teleportation (action = same state)

def sarsa(Q, model = model, alpha = 0.1, eps = 0.1, n_iter = 100):
    states = model.states
    terminals = model.terminals
    rewards = model.rewards
    gamma = model.gamma
    # random state (not terminal)
    state = np.random.choice(np.setdiff1d(np.arange(len(states)), terminals))
    # random action
    action = np.random.choice(Q[state].indices)
    new_state = action
    for t in range(n_iter):
        state_prev = state
        action_prev = action
        state = new_state
```

```

    if state in terminals: # if the new action gives terminal state
        # we start again from a new random state
        state = np.random.choice(np.setdiff1d(np.arange(len(model.states)),
                                                terminals))

        action = np.random.choice(Q[state].indices)
        Q[state_prev, action_prev] =
        (1 - alpha) * Q[state_prev, action_prev] + alpha * rewards[action_prev]
        # we removed the second term since state_prev was the last step
    else:
        # we choose the best action that has the highest expected value-action
        best_action = Q[state].indices[np.argmax(Q[state].data)]
        if np.random.random() < eps:
            # with probability of epsilon, we choose a random action
            action = np.random.choice(Q[state].indices)
        else:
            action = best_action
        Q[state_prev, action_prev] =
        (1 - alpha) * Q[state_prev, action_prev]
        + alpha * (rewards[action_prev] + gamma * Q[state, action])
    new_state = action
    return Q

```

Q-learning

This algorithm is also based on ϵ -greedy algorithm.

$$\pi(s) \leftarrow \begin{cases} a^* \in \operatorname{argmax}_a Q(s, a) & \text{with probability } 1 - \epsilon \\ \text{random} & \text{with probability } \epsilon \end{cases} \quad (4)$$

Estimation: unlike SARSA, Q-learning aims at updating the estimator using the best action at each iteration:

$$\forall t, Q(s_t, a_t) \stackrel{\alpha}{\leftarrow} r_t + \gamma \max_a Q(s_{t+1}, a)$$

The only modification from SARSA is the following line:

Listing 6: Q-Learning

```

Q[state_prev, action_prev] = (1 - alpha) * Q[state_prev, action_prev]
+ alpha * (rewards[action_prev] + gamma * np.max(Q[state].data))

```
