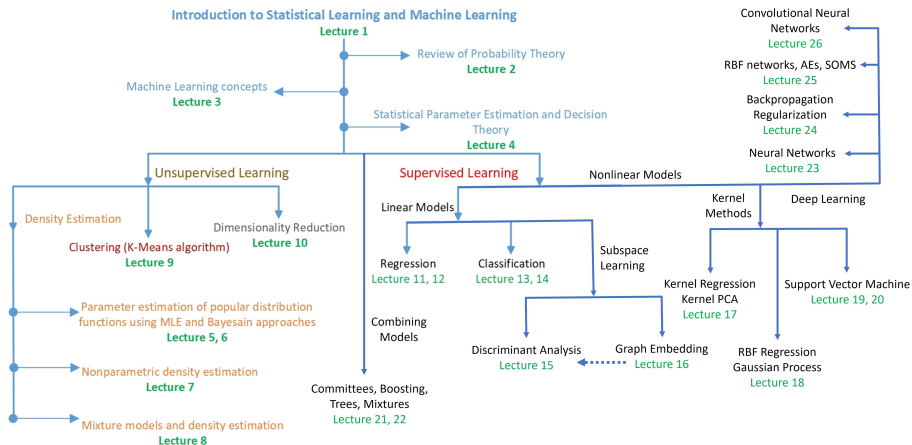


Statistical Learning and Machine Learning

Lecture 22 - Combining Models 2

October 13, 2021

Course overview and where do we stand



Why to combine multiple models?

Sometimes combining multiple models can lead to better performance compared to using only one:

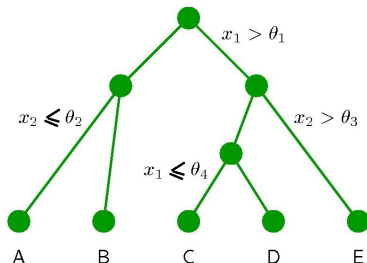
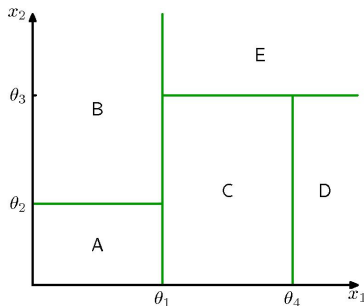
- *Committees*: Use L different models and then make predictions using the average of the predictions of these L models
- *Boosting*: Use multiple models *in a sequence* in which each model's training is depending on the models preceding it in the sequence
- *Decision trees*: Divide the space in multiple regions and train one model for each region.
- *Mixture of experts*: Train K models and combine them based on a probabilistic mixture of the form:

$$p(t|x) = \sum_{k=1}^K p(k|x)p(t|x, k) \quad (1)$$

Tree-based models

Process:

- 1 Partition the input space into cuboid regions
- 2 Assign a simple model (e.g. a constant) to each region



Tree-based models

Given $\{x_1, \dots, x_N\}$ and the corresponding targets $\{t_1, \dots, t_N\}$:

- ① Start with a single root node, corresponding to the whole input space
- ② Grow the tree by adding nodes one at a time:
 - Add a pair of leaf nodes to the existing tree corresponding to a candidate variable that is split in two
 - Determine the corresponding threshold for this variable
 - For a given choice of split variable and threshold, the optimal choice of predictive variable is given by the local average of the data
- ③ Repeated for all possible regions to be split, and the one that gives the smallest residual sum-of-squares error is retained.

Grow a large tree, using a stopping criterion based on the number of data points associated with the leaf nodes, and then prune back the resulting tree.

Tree-based models

Denote by T_0 the starting tree formed by $\tau = 1, \dots, |T|$ leaf nodes, with leaf node τ representing a region \mathcal{R}_τ including N_τ data points.

The optimal prediction for \mathcal{R}_τ is:

$$y_\tau = \frac{1}{N_\tau} \sum_{\mathbf{x}_n \in \mathcal{R}_\tau} t_n \quad (2)$$

For a regression problem, the contribution of \mathcal{R}_τ to the residual sum-of-squares is:

$$Q_\tau(T) = \sum_{\mathbf{x}_n \in \mathcal{R}_\tau} (t_n - y_\tau)^2 \quad (3)$$

and the pruning criterion is:

$$C(T) = \sum_{\tau=1}^{|T|} Q_\tau(T) + \lambda |T| \quad (4)$$

For classification problems $Q_\tau(T)$ is replaced by a classification loss (e.g. the cross-entropy loss).

Mixture models: Linear regression

We consider K linear models mapping ϕ vector to variable t :

- each having its own weight parameter vector w_k , $k = 1, \dots, K$
- all having the same noise variance, expressed by the precision parameter β
- The mixture distribution is:

$$p(t|\theta) = \sum_{k=1}^K \pi_k \mathcal{N}(t|w_k^T \phi, \beta^{-1}) \quad (5)$$

where π_k , $k = 1, \dots, K$ are the mixing coefficients and θ expresses the parameters of the mixture model (w_k , π_k , $k = 1, \dots, K$ and β)

Mixture models: Linear regression

Given a set of data points and target values $\{\phi_n, t_n\}$, $n = 1, \dots, N$ the log-likelihood functions is:

$$\ln p(t|\theta) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}(t_n | \mathbf{w}_k^T \phi_n, \beta^{-1}) \right) \quad (6)$$

To optimize the log-likelihood function we introduce a set of 1-of- K variables $Z = \{z_n\}$ (with $z_{nk} \in \{0, 1\}$) indicating which component of the mixture is responsible for generating the data point:

$$\ln p(t, Z, |\theta) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \ln \left(\pi_k \mathcal{N}(t_n | \mathbf{w}_k^T \phi_n, \beta^{-1}) \right) \quad (7)$$

We use the EM algorithm to estimate $\{\mathbf{w}_k, \pi_k\}$, $k = 1, \dots, K$ and β

Mixture models: Linear regression

Expectation Maximization algorithm:

Start by using θ^{old}

E-step: Use θ^{old} to evaluate the posterior probabilities (or *responsibilities*)

$$\gamma_{nk} = \mathbb{E}[z_{nk}] = p(k|\phi_n, \theta^{old}) = \frac{\pi_k \mathcal{N}(t_n | \mathbf{w}_k^T \phi_n, \beta^{-1})}{\sum_j \pi_j \mathcal{N}(t_n | \mathbf{w}_j^T \phi_n, \beta^{-1})} \quad (8)$$

$$Q(\theta, \theta^{old}) = \mathbb{E}[\ln p(\mathbf{t}, \mathbf{Z} | \theta)] = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left(\ln \pi_k + \ln \mathcal{N}(t_n | \mathbf{w}_k^T \phi_n, \beta^{-1}) \right) \quad (9)$$

Mixture models: Linear regression

Expectation Maximization algorithm:

M-step: Maximize $Q(\theta, \theta^{old})$ w.r.t. $\{w_k, \pi_k\}$, $k = 1, \dots, K$ and β keeping γ_{nk} fixed, and considering that $\sum_k \pi_k = 1$:

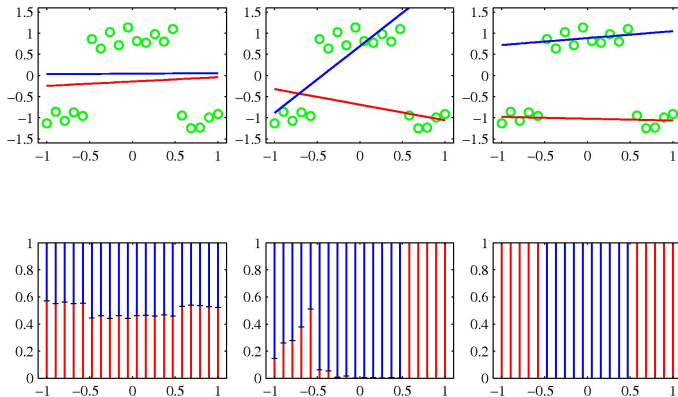
$$\pi_k = \frac{1}{N} \sum_{n=1}^N \gamma_{nk} \quad (10)$$

$$w_k = \left(\Phi^T R_k \Phi \right)^{-1} \Phi^T R_k t \quad (11)$$

$$\frac{1}{\beta} = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (t_n - w_k^T \phi_n)^2 \quad (12)$$

where $R_k = \text{diag}(\gamma_{nk})$: each data point ϕ_n is contributing to the calculation of the parameters of the k -th model according to the corresponding responsibility value γ_{nk} .

Mixture models: Linear regression

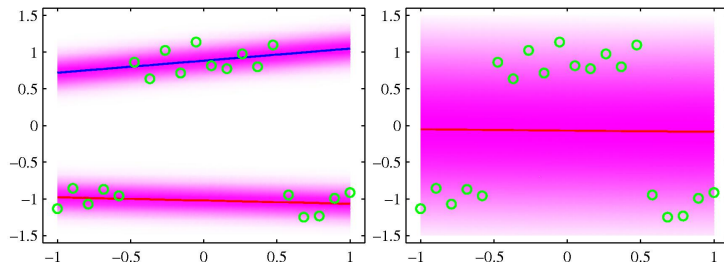


Regression using mixture of two linear regression models

Top: After 30 (left) and 50 iteration (center), final result (right)

Bottom: The responsibility values $p(k|\phi_n, \theta)$, $k = 1, 2$

Mixture models: Linear regression



Left: Predictive distribution for a mixture of two linear regression models
Right: Predictive density for a single linear regression model

Mixture models: Logistic regression

For a mixture of K logistic regression models:

$$p(t|\phi, \theta) = \sum_{k=1}^K \pi_k y_k^t [1 - y_k]^{1-t} \quad (13)$$

where:

- ϕ is the feature vector
- $y_k = \sigma(w_k^T \phi)$
- θ denotes the parameters $\{\pi_k, w_k\}$, $k = 1, \dots, K$

Mixture models: Logistic regression

Given a set of data points and target values $\{\phi_n, t_n\}$, $n = 1, \dots, N$ the likelihood functions is:

$$p(\mathbf{t}, \boldsymbol{\theta}) = \prod_{n=1}^N \left(\sum_{k=1}^K \pi_k y_{nk}^{t_n} [1 - y_{nk}]^{1-t_n} \right) \quad (14)$$

with $y_{nk} = \sigma(\mathbf{w}_k^T \phi_n)$ and $\mathbf{t} = [t_1, \dots, t_N]^T$.

To optimize the likelihood function we introduce a set 1-of- K variables $\mathbf{Z} = \{z_n\}$ indicating which component of the mixture is responsible for generating the data point:

$$p(\mathbf{t}, \mathbf{Z}, |\boldsymbol{\theta}) = \prod_{n=1}^N \prod_{k=1}^K \left(\pi_k y_{nk}^{t_n} [1 - y_{nk}]^{1-t_n} \right)^{z_{nk}} \quad (15)$$

We use the EM algorithm to estimate $\{\mathbf{w}_k, \pi_k\}$, $k = 1, \dots, K$ and β

Mixture models: Logistic regression

Expectation Maximization algorithm:

Start by using θ^{old}

E-step: Use θ^{old} to evaluate the posterior probabilities (or *responsibilities*)

$$\gamma_{nk} = \mathbb{E}[z_{nk}] = p(k|\phi_n, \theta^{old}) = \frac{\pi_k y_{nk}^{t_n} [1 - y_{nk}]^{1-t_n}}{\sum_j \pi_j y_{nj}^{t_n} [1 - y_{nj}]^{1-t_n}} \quad (16)$$

$$\begin{aligned} Q(\theta, \theta^{old}) &= \mathbb{E}_z[\ln p(t, Z|\theta)] \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left(\ln \pi_k + t_n \ln y_{nk} + (1 - t_n) \ln(1 - y_{nk}) \right) \end{aligned} \quad (17)$$

Mixture models: Logistic regression

Expectation Maximization algorithm:

M-step: Maximize $Q(\theta, \theta^{old})$ w.r.t. $\{w_k, \pi_k\}$, $k = 1, \dots, K$ keeping γ_{nk} fixed

$$\pi_k = \frac{1}{N} \sum_{n=1}^N \gamma_{nk} \quad (18)$$

$Q(\theta, \theta^{old})$ does not have a closed form solution for w_k and they need to be solved iteratively using the iterative reweighted least squares (IRLS) algorithm.

$$\nabla_k Q = \sum_{n=1}^N \gamma_{nk} (t_n - y_{nk}) \phi_n \quad (19)$$

$$H_k = -\nabla_k \nabla_k Q = \sum_{n=1}^N \gamma_{nk} y_{nk} (1 - y_{nk}) \phi_n \phi_n^T \quad (20)$$

Mixture models: Logistic regression

IRLS algorithm:

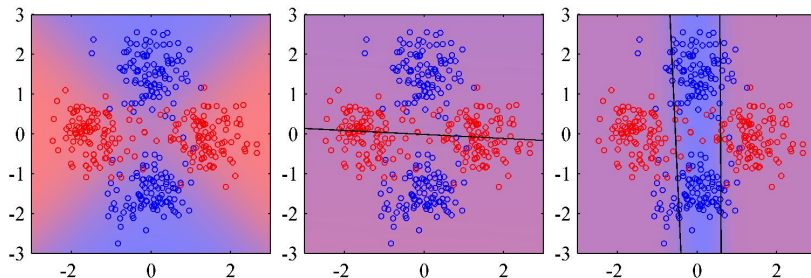
$$\nabla_k Q = \sum_{n=1}^N \gamma_{nk} (t_n - y_{nk}) \phi_n \quad (21)$$

$$H_k = -\nabla_k \nabla_k Q = \sum_{n=1}^N \gamma_{nk} y_{nk} (1 - y_{nk}) \phi_n \phi_n^T \quad (22)$$

where ∇_k denotes the derivative w.r.t. w_k

For fixed γ_{nk} the above equations are independent of $\{w_j\}$, $j \neq k$. Thus, we can solve for each w_k independently using the IRLS algorithm.

Mixture models: Logistic regression



Left: true probability of each class

Center: Result of fitting a single logistic regression model

Right: Result of fitting a mixture of two logistic regression models

Mixture of experts

We can increase the capability of mixture models by allowing the mixing coefficients be functions of the input variable too:

$$p(t|x) = \sum_{k=1}^K \pi_k(x) p_k(t|x). \quad (23)$$

The mixing coefficients $\pi_k(x)$ are known as *gating functions* and the component densities $p_k(t|x)$ are known as *experts*.

The gating functions need to satisfy the constraints:

$$0 \leq \pi_k(x) \leq 1 \quad \text{and} \quad \sum_{k=1}^K \pi_k(x) = 1. \quad (24)$$