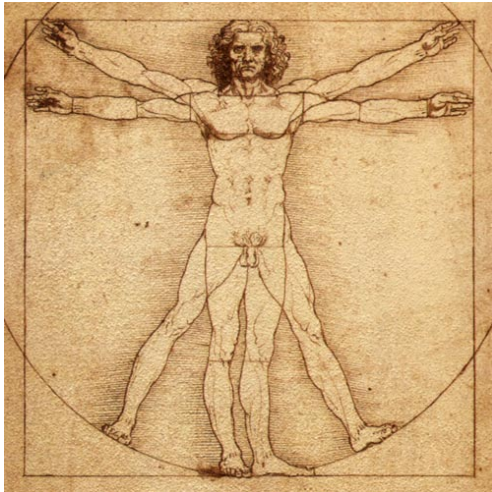


4F10: Ensemble Methods

Mark Gales

Michaelmas 2020

Renaissance Man

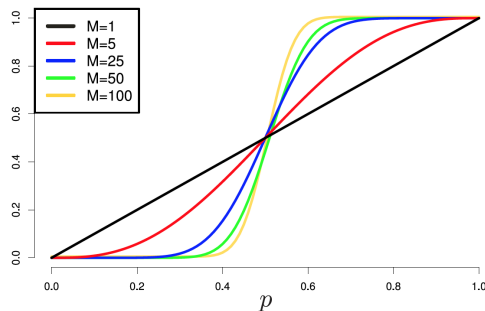


Binary Ensemble Classifier Majority Voting

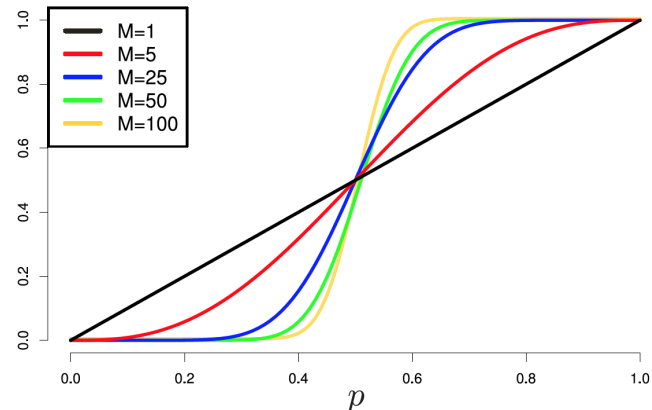
- Consider an ensemble of binary classifiers
 - each classifier's prediction is **independent** of other classifiers
 - select class with the most "votes" (majority voting)

$$P(\text{error}) = \sum_{i=\frac{M+1}{2}}^M \binom{M}{i} p^i (1-p)^{M-i}$$

- p is the probability of error from a single classifier



Ideal Binary Ensemble Error



- In practice not possible to ensure predictions independent
 - often related training data/criterion/topology

- Now consider an **ensemble of discriminative classifiers**
 - vary over topologies/parameters/training data etc.

$$\hat{\omega} = \arg \max_{\omega} \left\{ \sum_{\{\mathcal{M}\}} \int P(\omega | \mathbf{x}^*, \theta, \mathcal{M}) p(\theta | \mathcal{M}, \mathcal{D}) P(\mathcal{M} | \mathcal{D}) d\theta \right\}$$

- \mathcal{M} specifies an instance of a “model” topology
- θ are parameters of a particular model topology
- highly complicated (non parametric) distribution

- Now consider an **ensemble of discriminative classifiers**
 - vary over topologies/parameters/training data etc.

$$\hat{\omega} = \arg \max_{\omega} \left\{ \sum_{\{\mathcal{M}\}} \int P(\omega | \mathbf{x}^*, \theta, \mathcal{M}) p(\theta | \mathcal{M}, \mathcal{D}) P(\mathcal{M} | \mathcal{D}) d\theta \right\}$$

- \mathcal{M} specifies an instance of a “model” topology
 - θ are parameters of a particular model topology
 - highly complicated (non parametric) distribution
- Use a Monte-Carlo approximation - an **ensemble**, \mathcal{E} ,

$$\begin{aligned} \hat{\omega} &= \arg \max_{\omega} \left\{ \frac{1}{M} \sum_{j=1}^M P(\omega | \mathbf{x}^*; \theta^{(j)}) \right\} \\ \mathcal{E} &= \left\{ \theta^{(1)}, \dots, \theta^{(M)} \right\}; \quad \theta^{(i)} \sim p(\theta, \mathcal{M} | \mathcal{D}); \end{aligned}$$

- Posterior distribution, $p(\theta, \mathcal{M}|\mathcal{D})$, can be very complicated
 - distribution over model topologies
 - distribution over model parameters
 - distribution over model features/targets

- Posterior distribution, $p(\theta, \mathcal{M}|\mathcal{D})$, can be very complicated
 - distribution over model topologies
 - distribution over model parameters
 - distribution over model features/targets
- How to sample over posterior? Two distinct elements
 - **topology**: network/layers/activation functions
 - **parameters** for a topology distribution over the parameters
- Challenging to sample (efficiently) over topologies
- Range of practical approaches to generating **ensembles**
 - often limited links to underlying theoretical sampling

- Simple process based on partitioning the data:

1. from training data \mathcal{D} select subset $\tilde{\mathcal{D}}$ of size \tilde{n}

- selection can be with **replacement**

2. train model on $\tilde{\mathcal{D}}$

- for CML training with subset $\tilde{\mathcal{D}} = \{\{y_j, \mathbf{x}_j\}\}_{j=1}^{\tilde{n}}$

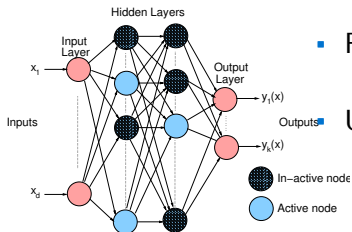
$$\theta^{(i)} = \arg \max_{\theta} \left\{ \sum_{j=1}^{\tilde{n}} \log(P(y_j | \mathbf{x}_j; \theta)) \right\}$$

- each “sample” at a “local optimal”

3. repeat until ensemble size generated

- Ensemble members built on a random sub-set of the data

Monte Carlo Dropout [4]



- Randomly de-activate nodes during training
 - improved regularisation of network
- Usually impact of dropout “undone”
 - “deweight” weights related to dropout rate

- Ensemble generated by de-activating random nodes
 - each random selection related to draw from “local posterior”
 - topology a subset of the original topology
- Allows ensemble to trained with only a single network

- For a given topology network parameters initialised used

$$\tilde{\theta}^{(i)} \sim p(\theta|\mathcal{M})$$

- $p(\theta|\mathcal{M})$ is the **prior** over the parameters for a topology
 - $\tilde{\theta}^{(i)}$ is used as model initialisation
- Each member of the ensemble then be trained
 - assume CML training with supervised data $\mathcal{D} = \{\{y_j, \mathbf{x}_j\}\}_{j=1}^n$

$$\theta^{(i)} = \arg \max_{\theta} \left\{ \sum_{j=1}^n \log(P(y_j|\mathbf{x}_j, \tilde{\theta}^{(i)}; \theta)) \right\}$$

- each “sample” at a “local optimal”

- Simple approach: select different configurations:
 1. number/size of layers
 2. nature of activation function/training cost function
 3. nature of targets
- Manually “sampling” from the “configuration posterior”
 - standard approaches used in many systems
 - each system at a local optimum
 - allows rich diversity of models to be used
 - expert knowledge to make systems complementary

- Simplest approach: average class posteriors from the ensemble

$$P(\omega|\mathbf{x}^*; \mathcal{E}) = \frac{1}{M} \sum_{i=1}^M P(\omega|\mathbf{x}^*; \boldsymbol{\theta}^{(i)})$$

- requires all models in the ensemble to be evaluated
- requires all model in the ensemble to be stored
- no explicit measure of structure in the selection

- Simplest approach: average class posteriors from the ensemble

$$P(\omega|\mathbf{x}^*; \mathcal{E}) = \frac{1}{M} \sum_{i=1}^M P(\omega|\mathbf{x}^*; \boldsymbol{\theta}^{(i)})$$

- requires all models in the ensemble to be evaluated
 - requires all model in the ensemble to be stored
 - no explicit measure of structure in the selection
- Interesting to consider:
 - can the ensemble be compressed to a single model?

- Consider **cross entropy** training criterion, supervised training
 - training data $\mathcal{D} = \{\{\mathbf{x}_1, y_1\}, \dots, \{\mathbf{x}_n, y_n\}\}$, $y_i \in \{\omega_1, \dots, \omega_K\}$

$$\mathcal{F}_{\text{ce}} = - \sum_{i=1}^n \log(P(y_i | \mathbf{x}_i; \boldsymbol{\theta}_s)) = - \sum_{i=1}^n \sum_{\omega} \delta(y_i, \omega) \log(P(\omega | \mathbf{x}_i; \boldsymbol{\theta}_s))$$

- $\delta(y_i, \omega)$ a Kronecker delta function, sum over all classes

- Consider **cross entropy** training criterion, supervised training
 - training data $\mathcal{D} = \{\{\mathbf{x}_1, y_1\}, \dots, \{\mathbf{x}_n, y_n\}\}$, $y_i \in \{\omega_1, \dots, \omega_K\}$

$$\mathcal{F}_{\text{ce}} = - \sum_{i=1}^n \log(P(y_i | \mathbf{x}_i; \boldsymbol{\theta}_S)) = - \sum_{i=1}^n \sum_{\omega} \delta(y_i, \omega) \log(P(\omega | \mathbf{x}_i; \boldsymbol{\theta}_S))$$

- $\delta(y_i, \omega)$ a Kronecker delta function, sum over all classes
- Modify targets based on a **teacher**

$$\mathcal{F}_{\text{ts}} = - \sum_{i=1}^n \sum_{\omega} P(\omega | \mathbf{x}_i; \boldsymbol{\theta}_T) \log(P(\omega | \mathbf{x}_i; \boldsymbol{\theta}_S))$$

- “soft” targets rather than “hard” targets to train $\boldsymbol{\theta}_S$
 - targets obtained from teacher network, $\boldsymbol{\theta}_T$
- Also called **model distillation**

Teacher-Student Ensemble Combination

- Originally used for model compression
 - student network, θ_S , **simpler** than teacher network, θ_T
 - can also be used for ensemble combination
- Replace a single teacher by an ensemble, \mathcal{E} , now

$$\mathcal{F}_{ts} = - \sum_{i=1}^n \sum_{\omega} P(\omega | \mathbf{x}_i; \mathcal{E}) \log (P(\omega | \mathbf{x}_i; \theta_S))$$

- Various forms possible for targets, for example

$$P(\omega | \mathbf{x}_i; \mathcal{E}) = \frac{1}{M} \sum_{j=1}^M P(\omega | \mathbf{x}_i; \theta^{(j)})$$

- where there are M models in the ensemble \mathcal{E}

Bias Variance Trade-Off (Reference Only)

- Consider the variance of a trained regression model $\hat{f}(\mathbf{x})$
 - the **ideal** regression model is $f(\mathbf{x})$: $y = f(\mathbf{x}) + \epsilon$
 - supervised training examples: $\mathcal{D} = \{\{\mathbf{x}_1, y_1\}, \dots, \{\mathbf{x}_n, y_n\}\}$
 - use data to train regression model $\hat{f}(\mathbf{x})$

$$\begin{aligned}\mathbb{E}\{(\hat{f}(\mathbf{x}) - y)^2\} &= \mathbb{E}\{\epsilon^2\} - (\mathbb{E}\{\epsilon\})^2 + \mathbb{E}\{\hat{f}(\mathbf{x})^2\} - (\mathbb{E}\{\hat{f}(\mathbf{x})\})^2 \\ &\quad + (\mathbb{E}\{f(\mathbf{x}) - \hat{f}(\mathbf{x})\})^2 \\ &= \text{Var}(\epsilon) + \text{Var}(\hat{f}(\mathbf{x})) + (\text{Bias}(\hat{f}(\mathbf{x})))^2\end{aligned}$$

- expectation is over all sets of training data drawn from $p(\mathbf{x}, y)$

Bias Variance Trade-Off - Regression

- Consider the variance of a trained regression model $\hat{f}(\mathbf{x})$
 - the **ideal** regression model is $f(\mathbf{x})$: $y = f(\mathbf{x}) + \epsilon$
 - supervised training examples: $\mathcal{D} = \{\{\mathbf{x}_1, y_1\}, \dots, \{\mathbf{x}_n, y_n\}\}$
 - use data to train regression model $\hat{f}(\mathbf{x})$

$$\begin{aligned}\mathbb{E}\{(\hat{f}(\mathbf{x}) - y)^2\} &= \mathbb{E}\{\epsilon^2\} - (\mathbb{E}\{\epsilon\})^2 + \mathbb{E}\{\hat{f}(\mathbf{x})^2\} - (\mathbb{E}\{\hat{f}(\mathbf{x})\})^2 \\ &\quad + (\mathbb{E}\{f(\mathbf{x}) - \hat{f}(\mathbf{x})\})^2 \\ &= \text{Var}(\epsilon) + \text{Var}(\hat{f}(\mathbf{x})) + (\text{Bias}(\hat{f}(\mathbf{x})))^2\end{aligned}$$

- expectation is over all sets of training data drawn from $p(\mathbf{x}, y)$
- Another view of generalisation:
 - note: can't do anything about **noise** in the data
 - improved generalisation \rightarrow simpler model \rightarrow higher bias
 - improved data modelling \rightarrow complex model \rightarrow higher variance

Ensemble Bias Variance Trade-Off (Regression)

- Consider a set of training data sets: $\mathcal{D}_1, \dots, \mathcal{D}_N$
 - train regressor for each training set: $\hat{f}^{(1)}(\mathbf{x}), \dots, \hat{f}^{(N)}(\mathbf{x})$
 - ensemble prediction

$$\hat{f}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \hat{f}^{(i)}(\mathbf{x})$$

- The expected **prediction error** is given by (true target y)

$$\mathbb{E} \{ (\hat{f}(\mathbf{x}) - y)^2 \} = \text{Var}(\epsilon) + \frac{1}{N} \text{Var}(\hat{f}^{(i)}(\mathbf{x})) + \left(1 - \frac{1}{N}\right) \text{Cov}(\hat{f}^{(i)}(\mathbf{x})) + (\text{Bias}(\hat{f}^{(i)}(\mathbf{x})))^2$$

- $\text{Var}(\hat{f}^{(i)}(\mathbf{x}))$: ind. class. var. $\mathbb{E} \{ (\hat{f}^{(i)}(\mathbf{x}))^2 \} - \mathbb{E} \{ (\hat{f}^{(i)}(\mathbf{x})) \}^2$
- $\text{Bias}(\hat{f}^{(i)}(\mathbf{x}))$: bias of individual classifier $\mathbb{E} \{ f(\mathbf{x}) - \hat{f}^{(i)}(\mathbf{x}) \}$
- Cov: covariance between classifiers $\mathbb{E} \{ (f^{(i)}(\mathbf{x}) - f^{(j)}(\mathbf{x}))^2 \}$

Adaboost (Reference)

- General approach for generating/combining an ensemble
 - converts multiple **weak** learners to **strong** learners
 - form described here is **AdaBoost**

- Consider a 2-class classification problem - training data

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}, \quad y_i \in \{-1, 1\}$$

- Combine classifiers as a recursive linear combination

$$\overline{\mathcal{F}}_m(\mathbf{x}_i) = \overline{\mathcal{F}}_{m-1}(\mathbf{x}_i) + \alpha_m \mathcal{F}(\mathbf{x}_i; \boldsymbol{\theta}^{(m)})$$

- How to train classifier $\mathcal{F}(\mathbf{x}_i; \boldsymbol{\theta}^{(m)})$ how to get weight α_m ?

AdaBoost (1)

- Define total cost function to minimise $E(\theta^{(m)})$
 - for correct classification $y_i = \mathcal{F}(\mathbf{x}_i; \theta^{(m)})$, $y_i \mathcal{F}(\mathbf{x}_i; \theta^{(m)}) = 1$
 - cost function to find $\theta^{(m)}$ - minimise

$$\begin{aligned} E(\theta) &= \sum_{i=1}^n \exp(-y_i \overline{\mathcal{F}}_m(\mathbf{x}_i)) \\ &= \sum_{i=1}^n \mathcal{G}_{m-1}(\mathbf{x}_i) \exp(-y_i \alpha_m \mathcal{F}(\mathbf{x}_i; \theta)) \end{aligned}$$

- Re-express criterion based on correct/incorrect using $\mathcal{F}(\mathbf{x}_i; \theta)$
 - simple to show that

$$\begin{aligned} E(\theta) &= \sum_{i=1}^n \mathcal{G}_{m-1}(\mathbf{x}_i) \exp(-\alpha_m) \\ &\quad + \sum_{y_i \neq \mathcal{F}(\mathbf{x}_i; \theta)} \mathcal{G}_{m-1}(\mathbf{x}_i) (\exp(\alpha_m) - \exp(-\alpha_m)) \end{aligned}$$

AdaBoost (2)

- Following on from previous cost function expression
 - optimal additional binary classifier, $\theta^{(m)}$, minimises

$$\begin{aligned}\theta^{(m)} &= \arg \min_{\theta} \left\{ \sum_{y_i \neq \mathcal{F}(\mathbf{x}_i; \theta)} \mathcal{G}_{m-1}(\mathbf{x}_i) \right\} \\ &= \arg \min_{\theta} \left\{ \sum_{y_i \neq \mathcal{F}(\mathbf{x}_i; \theta)} \exp(-y_i \overline{\mathcal{F}}_{m-1}(\mathbf{x}_i)) \right\}\end{aligned}$$

- can show that optimal weights are

$$\alpha_m = \frac{1}{2} \log \left(\frac{\sum_{y_i = \mathcal{F}(\mathbf{x}_i; \theta^{(m)})} \mathcal{G}_{m-1}(\mathbf{x}_i)}{\sum_{y_i \neq \mathcal{F}(\mathbf{x}_i; \theta^{(m)})} \mathcal{G}_{m-1}(\mathbf{x}_i)} \right)$$

- Has been extended to multiple classes

- [1] L. Breiman, "Bagging predictors," *Mach. Learn.*, vol. 24, no. 2, pp. 123–140, 1996.
- [2] C. Bucila, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *Proceedings of the Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Philadelphia, PA, USA, August 20-23, 2006*, 2006, pp. 535–541.
- [3] R. E. Schapire, "A brief introduction to boosting," in *Proceedings of the 16th International Joint Conference on Artificial Intelligence - Volume 2*, ser. IJCAI'99. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1999, pp. 1401–1406. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1624312.1624417>
- [4] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.