

Research Journal

Vadim Smolyakov (vss@mit.edu)

May 17, 2018

Contents

1	Bayesian Non-parametrics	3
1.1	Gaussian Process	3
1.2	Dirichlet Process	6
1.3	Construction of the DP	7
1.3.1	Blackwell-MacQueen Urn Scheme	7
1.3.2	Chinese Restaurant Process	8
1.3.3	Stick-Breaking Construction	8
1.4	Hierarchical Dirichlet Process (HDP)	9
1.4.1	HDP Hidden Markov Models	11
1.5	Dependent Dirichlet Process (DDP)	12
1.6	Indian Buffet Process	13
1.7	Infinite Mixture of Experts	14
1.8	DPMM	15
1.9	Fitting a DPMM model	17
2	Variational Inference	18
2.1	Application: Topic Models	22
2.2	Application: Image Denoising in Ising Model	25
2.3	Stochastic Variational Inference	28
3	Optimization	32
3.1	Simulated Annealing	32
3.2	Bayesian Optimization	33
3.3	Active Learning	36
3.3.1	Query Strategies	36
3.4	Submodular Optimization	39
3.4.1	Greedy Maximization	40
4	Information Planning	43
4.1	Introduction	43
4.2	Naive Bayes	44
4.3	Supervised LDA	47
4.4	Deep Neural Network	49
4.5	Discussion	51

5	Misc	53
5.1	Sequential Monte Carlo	53
5.2	Information Theory	55
5.2.1	Entropy	55
5.2.2	KL divergence	57
5.2.3	Mutual Information	59
5.3	Imbalanced Learning	61
5.4	Ensemble Methods	62
5.4.1	Bagging	63
5.4.2	Boosting	64
5.4.3	Stacking	66

1 Bayesian Non-parametrics

1.1 Gaussian Process

Gaussian processes (GPs) define a prior over functions that can be updated to a posterior once we have observed data. In a supervised setting, the function gives a mapping between the data points x_i and the target value y_i : $y_i = f(x_i)$. Gaussian processes infer a distribution over functions given the data $p(f|x, y)$ and then use it to make predictions given new data. A GP assumes that the function is defined at a finite and arbitrary chosen set of points x_1, \dots, x_n , such that $p(f(x_1), \dots, f(x_n))$ is jointly Gaussian with mean $\mu(x)$ and covariance $\Sigma(x)$, where $\Sigma_{ij} = \kappa(x_i, x_j)$ and κ is a positive definite kernel function.

Supervised learning can be divided into regression (prediction of continuous quantities) and classification (prediction of discrete class labels). Consider a simple regression problem:

$$f(x) = x^T w \quad y = f(x) + \epsilon \quad \epsilon \sim N(0, \sigma_n^2) \quad (1)$$

Assuming independent and identically distributed noise, we can write down the likelihood function:

$$p(y|X, w) = \prod_{i=1}^n p(y_i|x_i, w) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left\{-\frac{(y_i - x_i^T w)^2}{2\sigma_n^2}\right\} \sim N(Xw, \sigma_n^2 I) \quad (2)$$

In Bayesian framework, we need to specify a prior over the parameters: $w \sim N(0, \Sigma_p)$. Writing only the terms of the likelihood and the prior which depend on the weights, we get:

$$p(w|X, y) \propto \exp\left\{-\frac{1}{2\sigma_n^2} \|y - Xw\|^2\right\} \exp\left\{-\frac{1}{2} w^T \Sigma_p^{-1} w\right\} \quad (3)$$

$$\propto \exp\left\{-\frac{1}{2}(w - \bar{w})\left(\frac{1}{\sigma_n^2} X X^T + \Sigma_p^{-1}\right)(w - \bar{w})\right\} \quad (4)$$

$$\sim N\left(\frac{1}{\sigma_n^2} A^{-1} X y, A^{-1}\right) \quad (5)$$

where $A = \sigma_n^{-2} X X^T + \Sigma_p^{-1}$. Thus, we have a closed form posterior distribution over the parameters w . To make predictions using this equation, we need to invert the matrix A . Assuming the observations are noiseless, we want to predict the function outputs $y_* = f(x_*)$. Consider the following joint GP distribution:

$$\begin{pmatrix} f \\ f_* \end{pmatrix} \sim N\left(\begin{pmatrix} \mu \\ \mu_* \end{pmatrix}, \begin{pmatrix} K & K_* \\ K_*^T & K_{**} \end{pmatrix}\right) \quad (6)$$

where $K = \kappa(X, X)$, $K_* = \kappa(X, X_*)$ and $K_{**} = \kappa(X_*, X_*)$. Using standard rules for conditioning Gaussians, the posterior has the following form:

$$p(f_*|X_*, X, f) \sim N(f_*|\mu_*, \Sigma_*) \quad (7)$$

$$\mu_* = \mu(X_*) + K_*^T K^{-1} (f - \mu(X)) \quad (8)$$

$$\Sigma_* = K_{**} - K_*^T K^{-1} K_* \quad (9)$$

Figure 1 shows three functions drawn at random from a GP prior (left) and GP posterior (right) after observing five data points in the case of noise-free observations. The shaded area corresponds to two times the standard deviation around the mean (95% confidence region). We

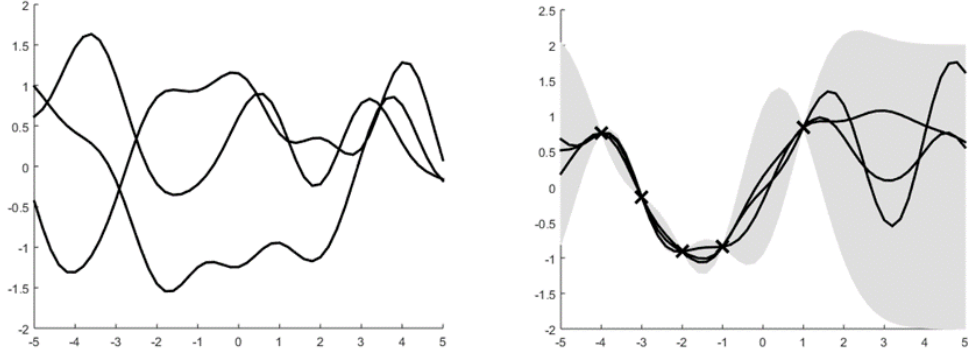


Figure 1: Gaussian process prior (left) and posterior (right).

can see that the model perfectly interpolates the training data and that the predictive uncertainty increases as we move further away from the observed data.

Since our algorithm is defined in terms of inner products in the input space, it can be lifted into feature space by replacing the inner products with $k(x, x')$, this is often referred to as the *kernel trick*. The kernel measures similarity between objects and it doesn't require pre-processing them into feature vector format. For a example, a common kernel function is a *radial basis function*:

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right) \quad (10)$$

In the case of a Gaussian kernel, the feature map lives in an infinite dimensional space. In this case, it is clearly infeasible to explicitly represent the feature vectors. Another commonly used kernel in Gaussian process regression is the *matern kernel*:

$$\kappa(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{l}\right)^\nu K_\nu\left(\frac{\sqrt{2\nu}r}{l}\right) \quad (11)$$

where $r = \|x - x'\|$, $\nu > 0$, $l > 0$, and K_ν is a modified Bessel function. As $\nu \rightarrow \infty$, this approaches the squared exponential kernel, if $\nu = \frac{1}{2}$, the kernel simplifies to $\kappa(r) = \exp(-r/l)$.

Assuming the observations are noisy, $y = f(x) + \epsilon$, where $\epsilon \sim N(0, \sigma_y^2)$, GP is not required to interpolate the data but rather it must be close to the observed data. The covariance of the observed noisy responses is $\text{cov}(y|X) = K + \sigma_y^2 I = K_y$, where we assumed that the noise terms were independently added to each observation. Hence, the posterior predictive density is:

$$p(f_*|X_*, X, y) \sim N(f_*|\mu_*, \Sigma_*) \quad (12)$$

$$\mu_* = K_*^T K_y^{-1} y \quad (13)$$

$$\Sigma_* = K_{**} - K_*^T K_y^{-1} K_* \quad (14)$$

In multiple dimensions, we can write down the *square exponential kernel* as follows:

$$\kappa_y(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2}(x_p - x_q)^T M (x_p - x_q)\right) + \sigma_y^2 \delta_{pq} \quad (15)$$

where $M = l^{-2}I$, l is the horizontal scale over which the function changes, σ_f^2 controls the vertical scale and σ_y^2 is the observation noise variance.

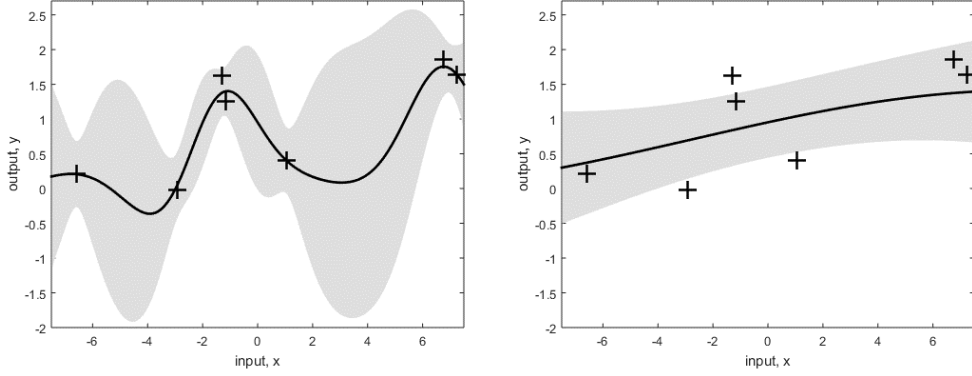


Figure 2: The effect of varying GP hyperparameters

Figure 2 shows the importance of choosing kernel hyperparameters and their impact on GP regression. We fix $\sigma_f^2 = 1$ and plot GP regression after 7 noisy observations. On the left, the length scale and noise variance were set to $(l, \sigma_y^2) = (1, 0.2)$. We can see that the mean function appears wiggly and has a small bias due to observation noise. On the right, the kernel parameters were set to $(l, \sigma_y^2) = (10, 0.8)$. The curve is much smoother due to large l and also shows higher observation noise.

We can use MCMC to learn GP hyperparameters by placing a prior on scale l , function variance σ_f^2 and noise variance σ_y^2 . We choose uninformative priors for GP hyper-parameters:

$$l \sim \text{Unif}(0, 10) \quad (16)$$

$$\sigma_f^2 \sim \exp\{\text{Unif}(-10, 5)\} \quad (17)$$

$$\sigma_y^2 \sim \exp\{\text{Unif}(-10, 5)\} \quad (18)$$

Using the No-U-Turn Sampler (NUTS) we obtain the following trace plots shown in Figure 3. We can see that posterior mode is close to ground truth hyper-parameters shown by solid red lines.

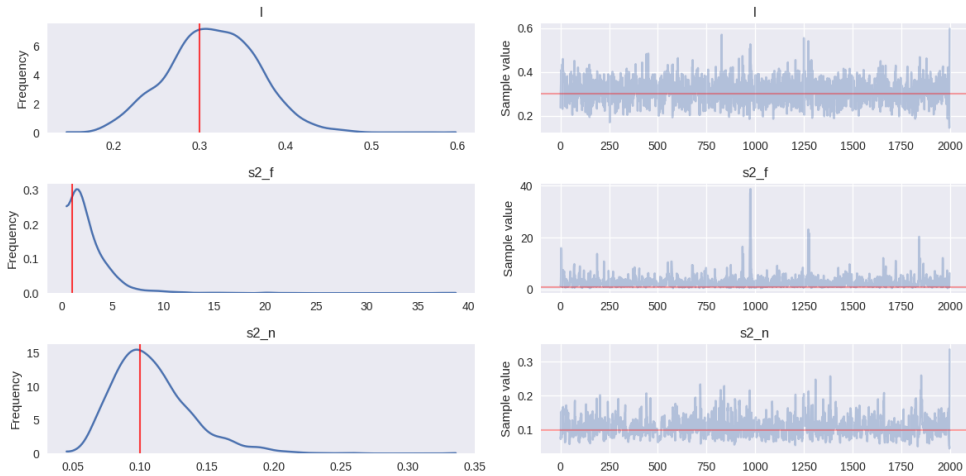


Figure 3: MCMC trace plots for GP hyperparameters.

Gaussian Processes (GPs) can be used for classification if the output of a GP is mapped

to the range $[0, 1]$. In the binary case, we define the model as $p(y_i|x_i) = \sigma(f(x_i))$, where we let $\sigma(z) = (1 + \exp(-z))^{-1}$. As with Bayesian logistic regression, the main difficulty in fitting GP classifier is that the Gaussian prior is not conjugate to the multinoulli likelihood. As a result several approximation algorithms can be used: Gaussian approximation, Expectation Propagation, variational and MCMC.

1.2 Dirichlet Process

The Dirichlet process is a stochastic process used in Bayesian non-parametric models. Each draw from a Dirichlet process is a discrete distribution. For a random distribution G to be distributed according to a DP, its finite dimensional marginal distributions have to be Dirichlet distributed. Let H be a distribution over Θ and α be a positive real number. We say that G is a Dirichlet process distributed with *base distribution* H and *concentration parameter* α if:

$$(G(A_1), \dots, G(A_r)) \sim \text{Dir}(\alpha H(A_1), \dots, \alpha H(A_r)) \quad (19)$$

for every finite measurable partition A_1, \dots, A_r of Θ , where Dirichlet distribution is defined as:

$$p(x_1, \dots, x_K) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K x_k^{\alpha_k - 1} \quad (20)$$

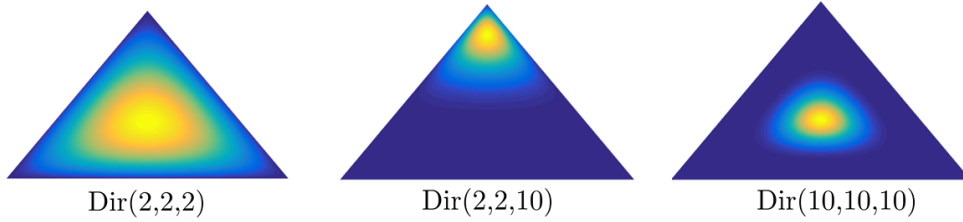


Figure 4: Dirichlet distribution: $\text{Dir}(\alpha_1, \alpha_2, \alpha_3)$ over the simplex.

The base distribution is the mean of the DP: $E[G(A)] = H(A)$, whereas the concentration parameter is the inverse variance: $V[G(A)] = H(A)(1 - H(A))/(\alpha + 1)$ for any measurable set $A \subset \Theta$. The larger the α , the smaller the variance, and the DP will concentrate more of its mass around the mean as shown in Figure 5.

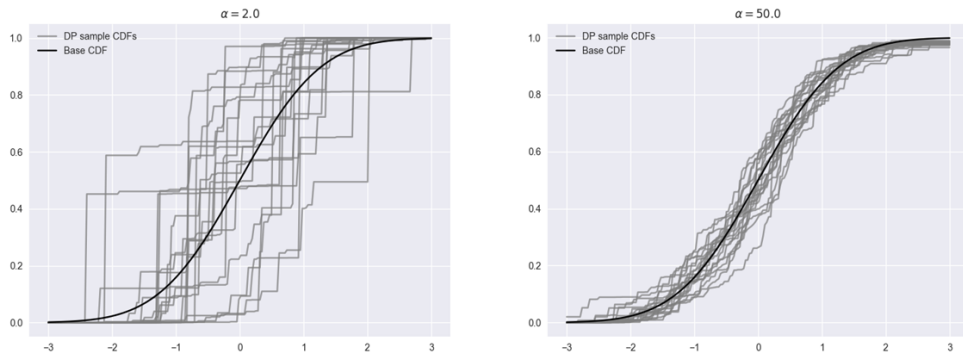


Figure 5: Base distribution CDF for $\alpha = 2$ and $\alpha = 50$.

Let $\theta_1, \dots, \theta_n$ be a sequence of independent draws from $G \sim DP(\alpha, H)$. We are interested in the posterior distribution of G given observed values of $\theta_1, \dots, \theta_n$. Let n_k be the number of

observed values in A_k , the by conjugacy the Dirichlet and multinomial distributions we have:

$$(G(A_1), \dots, G(A_r)) | \theta_1, \dots, \theta_n \sim \text{Dir}(\alpha H(A_1) + n_1, \dots, \alpha H(A_r) + n_r) \quad (21)$$

The posterior is also a DP with an updated concentration parameter and base distribution [40]:

$$G | \theta_1, \dots, \theta_n \sim DP(\alpha + n, \frac{\alpha}{\alpha + n} H + \frac{n}{\alpha + n} \frac{\sum_{i=1}^n \delta_{\theta_i}}{n}) \quad (22)$$

Note that the posterior base distribution is a weighted average between the prior base distribution H and the empirical distribution $\frac{1}{n} \sum_{i=1}^n \delta_{\theta_i}$. Therefore, α can be interpreted as the strength of the prior.

1.3 Construction of the DP

The Dirichlet Process can be represented in different ways and the representation influences the inference procedure.

1.3.1 Blackwell-MacQueen Urn Scheme

The Blackwell-MacQueen Urn Scheme, also commonly known as the Polya urn scheme, provides a way of constructing a sequence of draws $\theta_1, \theta_2, \dots, \theta_n \sim G$ where $G \sim DP(\alpha, H)$ without explicitly having to represent G . The posterior distribution with G marginalized out and $\theta_1 \sim H$ can be written as:

$$\theta_{n+1} | \theta_1, \dots, \theta_n \sim \frac{1}{\alpha + n} \left(\alpha H + \sum_{i=1}^n \delta_{\theta_i} \right) \quad (23)$$

The above equation describes a sequential process for drawing θ_i from G , which can be described by the following urn metaphor. First draw $\theta_1 \sim H$, paint a ball with that color and put it in the urn. Then, at each subsequent step n , either draw $\theta_n \sim H$ with probability $\alpha/(\alpha + n)$ and put a ball of that color in the urn, or with probability $n/(\alpha + n)$, randomly draw a ball from the urn, set θ_n to its color, then paint a new ball in that color and return both balls to the urn.

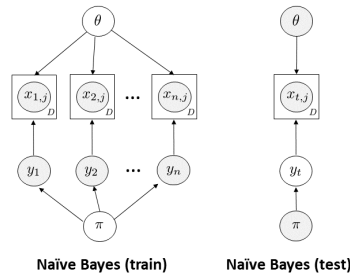


Figure 6: Naive Bayes graphical model

The Blackwell-MacQueen urn model has been used to show the existence of the DP. Starting from (23) we can construct a distribution over draws as follows:

$$P(\theta_1, \dots, \theta_n) = \prod_{i=1}^n P(\theta_i | \theta_1, \dots, \theta_{i-1}) \quad (24)$$

This random sequence is infinitely exchangeable: the probability of generating $\theta_1, \dots, \theta_n$ in that order is equal to the probability of generating them in any alternative order. Thus, for a given permutation σ , we have

$$P(\theta_1, \dots, \theta_n) = P(\theta_{\sigma(1)}, \dots, \theta_{\sigma(n)}) \quad (25)$$

The strength of infinite exchangeability lies in the following theorem:

Theorem 1.1 (*De Finetti, 1930s*) *A sequence of random variables x_1, x_2, \dots, x_n is infinitely exchangeable iff for all n :*

$$p(x_1, x_2, \dots, x_n) = \int p(\theta) \prod_{i=1}^n p(x_i|\theta) d\theta \quad (26)$$

Therefore, if we have exchangeable data, there must exist a parameter θ , a likelihood $p(x|\theta)$ and a distribution P on θ such that conditioned on θ the observations are independent as shown in Figure 6. In our setting, the prior over the random distribution $p(\theta)$ is the Dirichlet Process $DP(\alpha, H)$, thus establishing existence.

1.3.2 Chinese Restaurant Process

Chinese Restaurant Process (CRP) is a discrete-time stochastic process that is based on the clustering property of a DP. Let $\theta_1^*, \dots, \theta_m^*$ be unique values of draws $\theta_1, \dots, \theta_n$ and n_k be the number of repetitions of θ_k^* . Then, the predictive distribution can be written as:

$$\theta_{n+1}|\theta_1, \dots, \theta_n \sim \frac{1}{\alpha + n} \left(\alpha H + \sum_{k=1}^m n_k \delta_{\theta_k^*} \right) \quad (27)$$

Notice that θ_k^* will be repeated in proportion to n_k , the number of times it has already been observed. This is a *rich-gets-richer* phenomenon: the larger the cluster, the faster it will grow. The Chinese Restaurant Process takes its name from a metaphor describing its construction: imagine a Chinese restaurant which has an infinite number of tables, each of which can accommodate an infinite number of clusters. The first customer enters the restaurant and sits at the first table. The $n + 1$ st customer either joins an already occupied table k with probability proportional to n_k of customers already there or sits at a new table with probability proportional to α . Representing customers with integers and tables with clusters, CRP defines a partition on $[n]$. The number of occupied tables K approaches $\alpha \log(n)$ as n approaches infinity. And therefore, DP is a non-parametric prior that favors models whose complexity grows with the amount of data.

1.3.3 Stick-Breaking Construction

The stick breaking construction [34] represents draws $G \sim DP(\alpha, H)$ as a weighted sum of atoms (point masses). It is given as follows:

$$\beta_k \sim \text{Beta}(1, \alpha) \quad \theta_k^* \sim H \quad (28)$$

$$\pi_k = \beta_k \prod_{l=1}^{k-1} (1 - \beta_l) \quad G = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k^*} \quad (29)$$

This construction guarantees that $G \sim DP(\alpha, H)$. The name stick-breaking comes from the fact that one can interpret the π_k as the lengths broken off a unit length stick as shown in Figure 7 for $\alpha = 5$ and $\alpha = 10$.

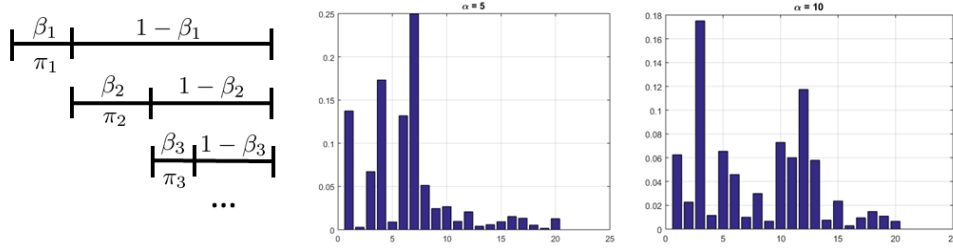


Figure 7: Stick-Breaking Weights

1.4 Hierarchical Dirichlet Process (HDP)

Hierarchical Dirichlet Processes model problems involving groups of data, where each observation within a group is a draw from a mixture model and mixture components are shared between groups. In each group, the number of components is learned from data using a Dirichlet Process prior. In addition, the base measure for the child Dirichlet processes is itself distributed according to a Dirichlet process. Since a draw from the global DP is a discrete distribution, the group-level DPs share mixture components. The HDP model can be summarized as follows and is shown in Figure 8:

$$G_0 | \gamma, H \sim DP(\gamma, H) \quad (30)$$

$$G_j | \alpha, G_0 \sim DP(\alpha, G_0) \quad (31)$$

The distribution G_0 varies around the prior H , with the amount of variability given by γ . The

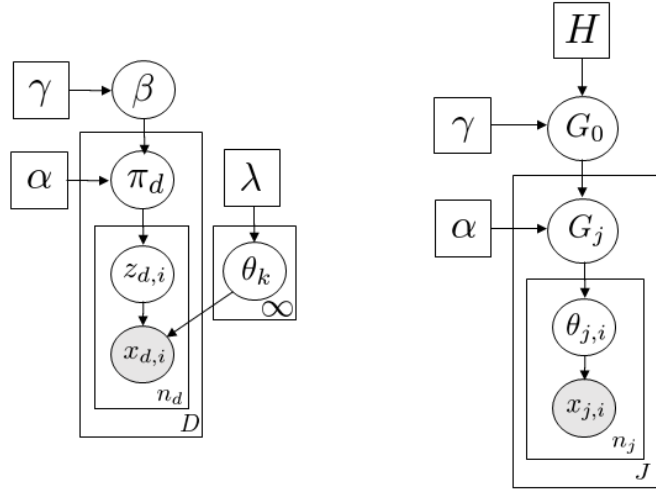


Figure 8: HDP Graphical Model

actual distribution G_j over $\theta_{j,i}$ in j^{th} group deviates from G_0 with the amount of variability given by α . If we expect the amount of variability between groups to be different, we can use a separate concentration parameter α_j for each group j .

Given the global DP prior G_0 , we can express it using a stick-breaking representation: $G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\phi_k}$, where $\phi_k \sim H$ and $\beta_k \sim \text{GEM}(\gamma)$. Since G_0 has support at the points ϕ_k , each G_j has support at these points as well and can be written as: $G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta_{\phi_k}$. The group weights π_j are related to global weights β as follows [39]:

$$\beta'_k \sim \text{Beta}(1, \gamma) \quad \beta_k = \beta'_k \prod_{l=1}^{k-1} (1 - \beta'_l) \quad (32)$$

$$\pi'_{jk} \sim \text{Beta}(\alpha\beta_k, \alpha(1 - \sum_{l=1}^k \beta_l)) \quad \pi_{jk} = \pi'_{jk} \prod_{l=1}^{k-1} (1 - \pi'_{jl}) \quad (33)$$

An analog of the Chinese restaurant process for HDPs is the Chinese restaurant franchise (CRF): the metaphor is extended to allow multiple restaurants which share a set of dishes. In this interpretation, each group defines a separate restaurant in which customers (observations) x_{ji} sit at tables (clusters) t_{ji} . Each table shares a single dish (parameter) θ_{ji} , which is ordered from a menu G_0 shared among restaurants (groups) as shown in Figure 9



Figure 9: Chinese Restaurant Franchise (CRF). Each restaurant is represented by a rectangle. Customers x_{ji} are seated at tables. At each table a dish ϕ_k is served from a global menu.

Let ϕ_1, \dots, ϕ_K denote the K atoms distributed according to H : this is the global menu of dishes. Let ψ_{jt} represent table-specific choice of dishes; in particular ψ_{jt} is the dish served at table t in restaurant j . Note that each θ_{ji} is associated with one ψ_{jt} and each ψ_{jt} is associated with one ϕ_k . In the CRF metaphor, customer i in restaurant j sat at table t_{ji} , while table t in restaurant j served dish k_{jt} .

The number of tables in restaurant j serving dish k is denoted m_{jk} , and the number of customers in restaurant j at table t eating dish k is n_{jtk} . Marginal counts are represented with dots. For example, $n_{j\cdot}$ and $m_{j\cdot}$ represent the number of customers and tables, respectively, in restaurant j .

We can compute the marginals under HDP when G_0 and G_j are integrated out using (22):

$$\theta_{ji} | \theta_{j1}, \dots, \theta_{jn}, \alpha, G_0 \sim \sum_{t=1}^{m_{j\cdot}} \frac{n_{jti}}{n + \alpha} \delta_{\psi_{jt}} + \frac{\alpha}{n + \alpha} G_0 \quad (34)$$

If a term in the first summation is chosen, we set $\theta_{ji} = \psi_{jt}$. If the second term is chosen then we increment $m_{j\cdot}$ by one and draw $\psi_{jm_{j\cdot}} \sim G_0$. To integrate out G_0 , we can use (22) again:

$$\psi_{jt} | \psi_{11}, \psi_{12}, \dots, \psi_{21}, \dots, \gamma, H \sim \sum_{k=1}^K \frac{m_{\cdot k}}{m_{\cdot\cdot} + \gamma} \delta_{\phi_k} + \frac{\gamma}{m_{\cdot\cdot} + \gamma} H \quad (35)$$

To summarize, for each j and i , we first sample θ_{ji} using (34). If a new sample from G_0 is needed, we use (35) to obtain a new sample ψ_{jt} .

An on-line variational inference algorithm for Hierarchical Dirichlet Process [43] was used to fit a topic model on the 20newsgroups dataset. The dataset consists of 11,314 documents and over 100K unique tokens. Standard text pre-processing was used including tokenization, stop-word removal and stemming. A compressed dictionary of $4K$ words was constructed by

filtering out tokens that appear in less than 5 documents and more than 50% of the corpus. The top-level truncation was set to $T = 20$ topics and the second level truncation was set to $K = 8$ topics. The concentration parameters were chosen as $\gamma = 1.0$ at the top-level and $\alpha = 0.1$ at the group level to yield a broad range of shared topics that are concentrated at the group level. Figure 10 shows a sample of the global level topics inferred by online variational HDP algorithm. We can find topics about autos, politics and for sale items that correspond to the target labels of the 20newsgroups dataset.

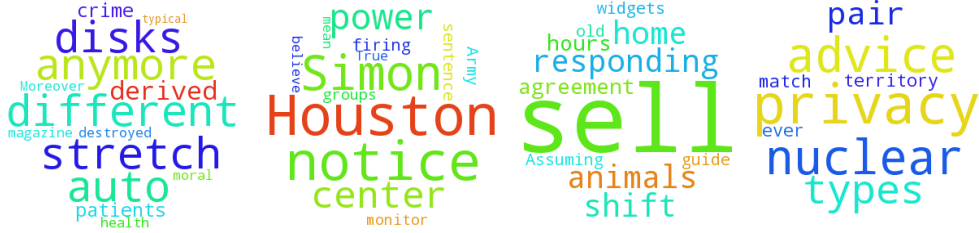


Figure 10: Sample of HDP topics inferred using online variational bayes algorithm on 20news-groups dataset.

1.4.1 HDP Hidden Markov Models

The Hierarchical Dirichlet Process (HDP) can be used to define a prior distribution on transition matrices over countably infinite state spaces. The HDP-HMM is known as an infinite Hidden Markov Model where the number of states is inferred automatically. To consider a non-parametric variant of the HMM, we must consider a set of DPs, one for each value of the current state. In addition, the DPs must be linked because we want the same set of next states to be reachable from each of the current states. This relates directly to HDP, where the atoms associated with state-conditional DPs are shared.

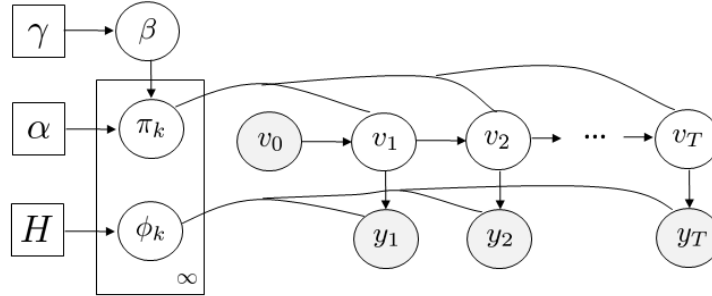


Figure 11: HDP-HMM graphical model

We can describe the HDP-HMM model in Figure 11 using the stick-breaking construction. The parameters have the following distribution:

$$\beta|\gamma \sim \text{GEM}(\gamma) \quad \pi_k|\alpha, \beta \sim \text{DP}(\alpha, \beta) \quad \phi_k|H \sim H \quad (36)$$

for each $k = 1, 2, \dots$ while for time steps $t = 1, \dots, T$ the state and observation distributions are:

$$v_t|v_{t-1}, \pi_k \sim \pi_{v_{t-1}} \quad y_t|v_t, \phi_k \sim F(\phi_{v_t}) \quad (37)$$

given a starting state v_0 . Each π_j is a DP draw and is interpreted as the transition distribution from state j . The π_j s are linked through DP draws parameterized by the same discrete measure β . Therefore, the states are shared between different DP draws.

Besides Gibbs sampling, one way to compute the posterior for infinite HMM is using a *beam sampling* algorithm [42]. Beam sampling combines two ideas: slice sampling and dynamic programming to sample whole state trajectories. Slice sampling is applied to limit the number of states to a finite number in each time step of the iHMM. The underlying idea is to limit the beam search to a small number of states so that a good trajectory can be found. Beam sampling is an MCMC method that guarantees convergence to the true posterior.

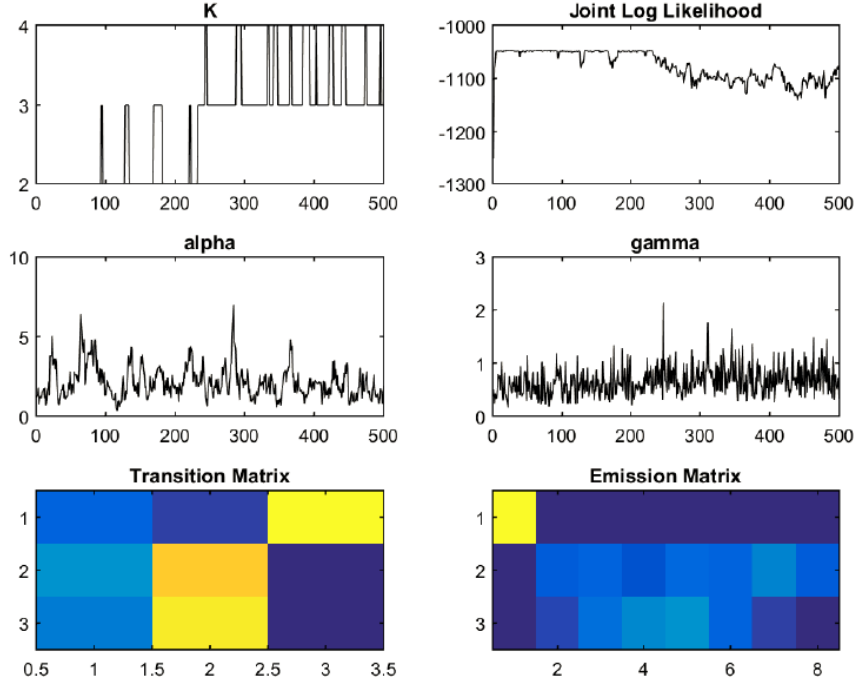


Figure 12: HDP-HMM inference results using the beam sampling algorithm [42].

Figure 12 shows posterior inference results for synthetically generated time-series data. The ground truth data was generated using 4 unique states for given transition and emission matrices. From the top-left plot in Figure 12, we can see that the beam sampler discovered three transition states at iteration 250. We can also see the samples of α and γ concentration parameters for the HDP-HMM as well as the transition and emission matrices.

1.5 Dependent Dirichlet Process (DDP)

The earlier part of Bayesian non-parametric literature focused on problems where a single distribution is assigned a non-parametric prior. However, in many applications, the objective is modelling a collection of distributions used in temporal and spatial processes. The Dirichlet process assumes that observations are exchangeable and therefore the data points have no inherent ordering that influences their labelling. This assumption is invalid for modelling temporal and spatial processes in which the order of data points plays a critical role in creating meaningful clusters.

The dependent Dirichlet process (DDP) originally formulated by MacEachern [24] provides a non-parametric prior over evolving mixture models. A construction of the DDP built on Poisson process [23] led to the development of the DDP mixture model (DDPMM) which generalizes DPMM by including birth, death and transition processes for the clusters in the model. In

addition, a low-variance approximations to DDPMM have been derived leading to a dynamic clustering algorithm [5].

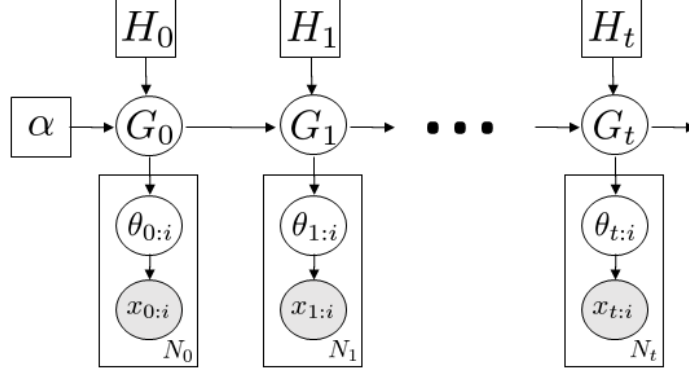


Figure 13: Dependent Dirichlet Process (DDP) graphical model.

Under time-varying setting, it is natural to introduce different DP priors for different time steps as shown in Figure 13. The generative model can be written as:

$$D_t \sim DP(\alpha, H_t) \quad (38)$$

$$\theta_{t:i}|D_t \sim D_t \quad \text{for } i = 1, \dots, n_t, \quad t = 0, \dots, T \quad (39)$$

$$x_{t:i}|\theta_{t:i} \sim F(\theta_{t:i}) \quad \text{for } i = 1, \dots, n_t, \quad t = 0, \dots, T \quad (40)$$

A Poisson-based construction of DDP [23] exploits the connection between Poisson and Dirichlet processes. In particular, by applying operations that preserve *complete randomness* to the underlying Poisson processes: superposition, subsampling and point transition, a new Poisson and therefore a new Dirichlet process is produced. Therefore, a Markov chain of Dirichlet processes can be written as:

$$D_t = T(S_q(D_{t-1})) \bigoplus H_t, \quad \text{where } H_t \sim DP(\alpha, H) \quad (41)$$

where S_q is an acceptance function, T is a probabilistic transition, combined with new terms from innovation process H_t to form D_t .

1.6 Indian Buffet Process

The Indian Buffet Process (IBP) is a stochastic process defining a probability distribution over sparse binary matrices with a finite number of rows and an infinite number of columns. This distribution is suitable to use as a prior for models with potentially infinite number of features. The form of the prior ensures that only a finite number of features will be present in any finite set of observations but more features may appear as more data points are observed.

Let Z be a $N \times K$ binary matrix indicating the presence or absence of a latent feature. The IBP places the following prior on Z [16]:

$$p(Z) = \frac{\alpha^K}{\prod_{i=1}^N K_1^{(i)}!} \exp\{-\alpha H_N\} \prod_{k=1}^K \frac{(N - m_k)!(m_k - 1)!}{N!} \quad (42)$$

where K is the number of nonzero columns in Z , m_k is the number of ones in column k of Z , H_N is the N^{th} harmonic number, and K_h is the number of occurrences of the non-zero binary vector h among the columns in Z . The parameter α controls the expected number of features present in each observation.

In the Indian Buffet Process (IBP), the rows of Z correspond to customers and the columns correspond to dishes in an infinitely long buffet. The first customer takes the first $\text{Poisson}(\alpha)$ dishes. The i th customer then takes dishes that have been previously sampled with probability m_k/i , where m_k is the number of people who have already sampled dish k . He also takes $\text{Poisson}(\alpha/i)$ new dishes. Then, z_{nk} is one if customer n tried k th dish and zero otherwise.

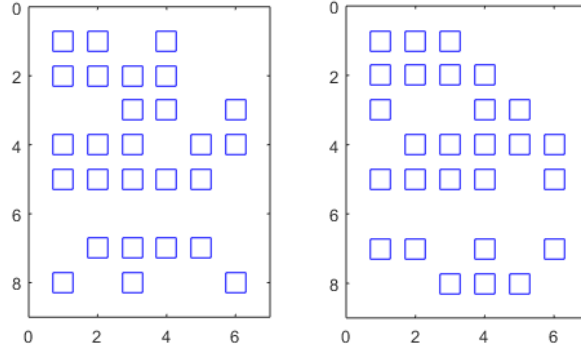


Figure 14: Original binary feature matrix (left) and its left-ordered form (right)

This process is infinitely exchangeable for an equivalence class of binary matrices defined by a left-ordered many-to-one function. $\text{lof}(Z)$ is obtained by ordering the columns of the binary matrix Z from left to right by the magnitude of the binary number expressed by that column, taking the first row as the most significant bit. The left ordering of a binary matrix is shown in Figure 14

In this section, we focus on variational inference procedures for the linear-Gaussian likelihood model [9]. Let X be a $N \times D$ matrix where each of the N rows contains a D -dimensional observation. We focus on a model where X can be approximated as:

$$X_{N \times D} = Z_{N \times K} \times A_{K \times D} + \epsilon \quad (43)$$

where Z is a binary feature matrix, the values for feature k are stored in row k of A , and ϵ is measurement noise. Figure 15 shows the predicted binary feature matrix Z and lower bound on marginal likelihood using variational inference algorithm for IBP [9]. The concentration parameter was set to $\alpha = 1$ and the feature truncation level was set to $K = 6$.

1.7 Infinite Mixture of Experts

Just as Dirichlet process mixtures can be thought of as infinite mixture models that select the number of active components as part of inference, infinite mixture of experts can be thought of selecting active components from an infinite number of regressions. This flexibility makes it a powerful tool for nonparametric Bayesian data analysis.

Infinite mixture of experts generalizes the Dirichlet Process Mixture Model by allowing the

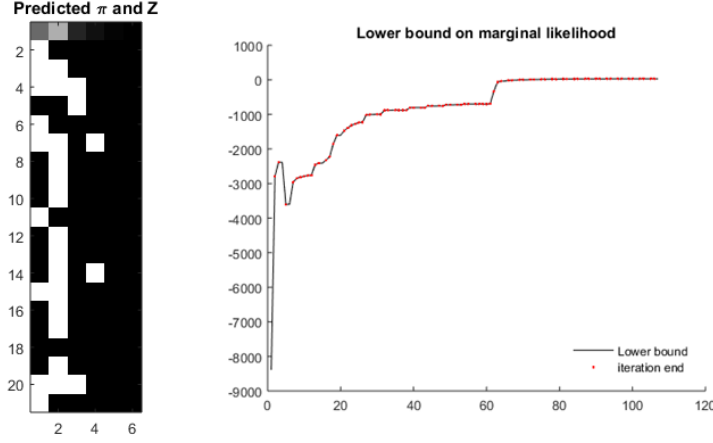


Figure 15: Predicted IBP binary feature matrix Z using variational inference.

mixture weights and component means to vary according to the input predictor x :

$$y|x \sim \sum_{i=1}^{\infty} w_i|x \times N(\mu_i|x, \tau_i^{-1}) \quad (44)$$

To determine the conditional mixture weights, we can use probit stick-breaking process:

$$v_i|x = \Phi(\alpha_i + \beta_i x) \quad (45)$$

$$w_i|x = v_i|x \times \prod_{j=1}^{i-1} (1 - v_j|x) \quad (46)$$

where Φ is the cumulative distribution function of the standard normal. To model multiple regressions, we define the mean of the mixture of experts as:

$$\mu_i|x \sim \gamma_i + \delta_i x \quad (47)$$

Figure 16 shows an infinite mixture of experts applied to LIDAR dataset that shows locally linear relationship between range and log ratio as well as heteroskedastic observation noise. The infinite mixture was truncated at $K = 20$ components which is reasonable considering that there are only three linear regressions having non-zero posterior expected weight.

Normal priors were set on regression coefficients for mixture mean and mixture weights, in addition $\tau_i \sim \text{Gamma}(1, 1)$ prior was set on the precision of the Gaussian mixture. The figure also shows a 95% credible interval that captures well the time-varying nature of observation noise. The flexibility of infinite mixture of experts and the ability to modularly specify the conditional mixture weights and component densities makes it a very useful nonparametric Bayesian model.

1.8 DPMM

The Dirichlet Process Mixture Models (DPMM) belong to a class of *infinite mixture models*, in which we do not impose any prior knowledge on the number of clusters K . DPMM models learn the number of clusters from the data using a non-parametric prior based on the Dirichlet Process (DP). Automatic model selection leads to computational savings of cross validating the model for multiple values of K .

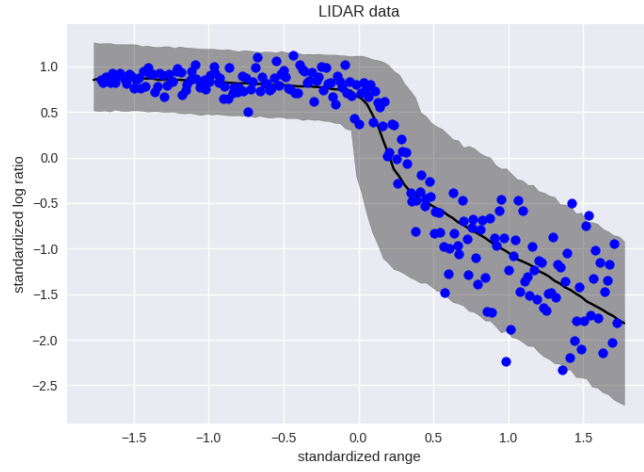


Figure 16: Infinite Mixture of Experts Regression.

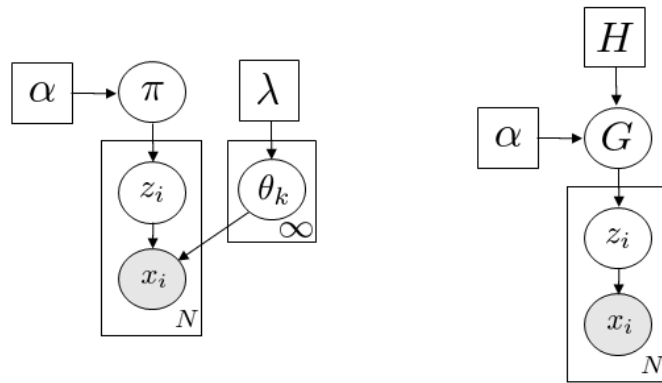


Figure 17: DPMM Graphical Model

Consider the graphical model of the DPMM in Figure 17. For each data point x_i , there's a corresponding label z_i that assigns the point to one of the clusters with mixture weight π_k and parameters $\theta_k = \{\mu_k, \Sigma_k\}$. The generative model can be written as follows:

$$p(x_i|z_i = k, \theta) = p(x_i|\theta_k) \quad (48)$$

$$p(z_i = k|\pi) = \pi_k \quad (49)$$

$$p(\pi|\alpha) = \text{Dir}(\pi|\alpha) \quad (50)$$

$$p(\theta_k|\beta) = \text{NIW}(\mu_k, \Sigma_k|m_0, \kappa_0, \nu_0, S_0) \quad (51)$$

An equivalent representation of this model can be written as:

$$G(\theta) = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k} \quad (52)$$

where $\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$ and $\theta_k \sim H$. Therefore, G is an infinite mixture of cluster functions of delta functions. In practice, we can construct G using a *stick-breaking construction*:

$$\beta_k \sim \text{Beta}(1, \alpha) \quad (53)$$

$$\pi_k = \beta_k \prod_{l=1}^{k-1} (1 - \beta_l) = \beta_k (1 - \sum_{l=1}^{k-1} \pi_l) \quad (54)$$

This is commonly denoted as $\pi \sim \text{GEM}(\alpha)$. The number of generated mixture components increases with α . The hyper-parameter α controls the expected number of clusters:

$$E[K] = \alpha \times \log(1 + n/\alpha) \quad (55)$$

$$\text{VAR}[K] = \alpha \times \log(1 + n/\alpha) \quad (56)$$

Thus, the number of clusters grows logarithmically with the number of data points and in direct proportion to α as shown in Figure 18

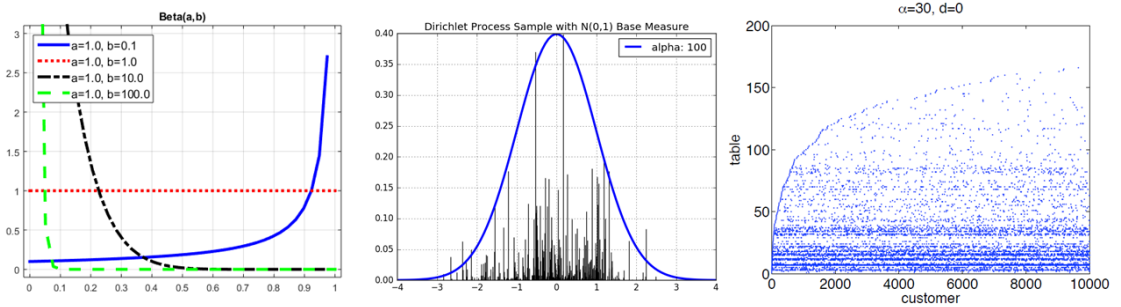


Figure 18: a) Beta(1, α) b) DP samples c) DP cluster growth

1.9 Fitting a DPMM model

One way to fit a DPMM is to modify the collapsed Gibbs sampler for a finite mixture model [26]:

$$p(z_i = k|z_{-i}, x, \alpha, \beta) \propto p(z_i = k|z_{-i}, \alpha) p(x_i|x_{-i}, z_i = k, z_{-i}, \beta) \quad (57)$$

The first term is given by the Chinese Restaurant Process (CRP):

$$p(z_i = k|z_{-i}, \alpha) = \begin{cases} \frac{N_k}{N + \alpha - 1} & \text{if } k \text{ is occupied} \\ \frac{\alpha}{N + \alpha - 1} & \text{if } k \text{ is a new cluster} \end{cases} \quad (58)$$

The second term is the posterior predictive and can be computed as follows:

$$p(x_i|x_{-i}, z_i = k, z_{-i}, \beta) = p(x_i|x_{k \setminus i}, \beta) = \frac{p(x_k|\beta)}{p(x_{k \setminus i}|\beta)} \quad (59)$$

We can compute both the numerator and the denominator above if we can find an expression for the marginal $p(x)$:

$$p(x) = \int_{\mu} \int_{\Sigma} p(x, \mu, \Sigma) d\mu d\Sigma = \int_{\mu} \int_{\Sigma} p(x|\mu, \Sigma) p(\mu, \Sigma|\beta) d\mu d\Sigma \quad (60)$$

$$= (2\pi)^{-ND/2} \frac{Z_{NIW}(D, \kappa_N, \nu_N, S_N)}{Z_{NIW}(D, \kappa_0, \nu_0, S_0)} = \pi^{-ND/2} \frac{\kappa_0^{D/2} |S_0|^{\nu_0/2}}{\kappa_N^{D/2} |S_N|^{\nu_N/2}} \prod_{i=1}^D \frac{\Gamma(\frac{\nu_N+1-i}{2})}{\Gamma(\frac{\nu_0+1-i}{2})} \quad (61)$$

Therefore, we can write the second term as follows:

$$\log p(x_i|x_{k \setminus i}) = z(D, N+1, \kappa_{N+1}, \nu_{N+1}, S_{N+1}) - z(D, N, \kappa_N, \nu_N, S_N) \quad (62)$$

$$z(D, N, \kappa, \nu, S) = -\frac{ND}{2} \log \pi - \frac{D}{2} \log \kappa - \frac{\nu}{2} \log |S| + \sum_{i=1}^D \log \Gamma\left(\frac{\nu+i-1}{2}\right) \quad (63)$$

We can summarize, the collapsed Gibbs sampler for an infinite Gaussian mixture model in Algorithm 1.

Algorithm 1 Collapsed Gibbs for DPMM

- 1: Initialize labels z
 - 2: for $t = 1, 2, \dots, T$ do
 - 3: for $i = 1, 2, \dots, N$ do
 - 4: remove x_i 's statistics from component z_i
 - 5: delete empty component
 - 6: for $k = 1, 2, \dots, K+1$ do
 - 7: compute $\log p(z_i = k|z_{-i}, x, \alpha, \beta) \sim$
 - 8: $\log p(z_i = k|z_{-i}, \alpha) + \log p(x_i|x_{k \setminus i}, \beta)$
 - 9: sample $z_i = k_{new}$ from $p(z_i = k|z_{-i}, x, \alpha, \beta)$
 - 10: create a new cluster if $z_i = K+1$
 - 11: add x_i 's statistics to component $z_i = k_{new}$
 - 12: end for
 - 13: end for
-

Figure 19 shows the clustering results of DPMM collapsed gibbs sampler with Gaussian base measure and $\alpha = 1$ after 100 iterations on a synthetic dataset of $1K$ points in 2D. The sampler was initialized with $K_{init} = 2$ clusters and correctly identified all 5 clusters in the data.

Figure 20 shows the clustering results of DPMM on categorical data. Documents from a reduced NIPS dataset with 300 docs, 5K vocab and 130K words were grouped into 11 clusters. The top words for the first 4 clusters are displayed in Figure 20. The results indicate meaningful division of articles by topics about neural networks, gaussian mixtures, computer vision and reinforcement learning. The gibbs sampler was initialized with $K_{init} = 2$ and $\alpha = 1$ and covered after 20 iterations through the corpus.

2 Variational Inference

This section focuses on a class of approximate inference algorithms based on variational inference. The basic idea is to choose an approximation $q(x)$ from a tractable family of distributions

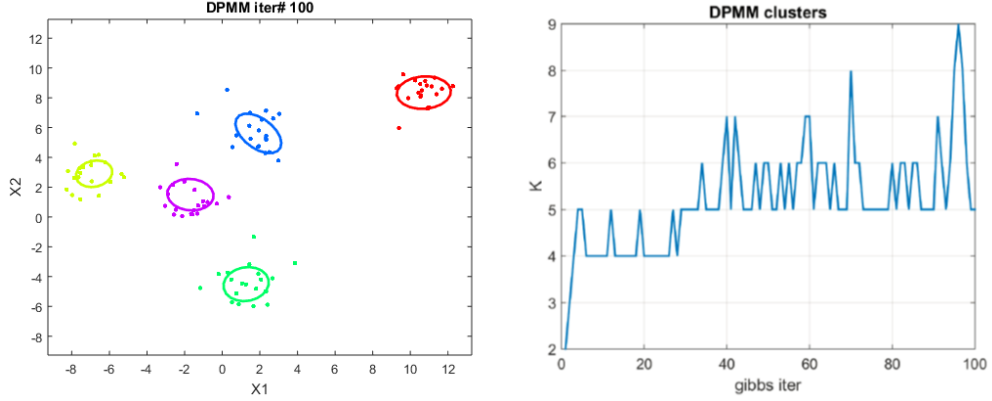


Figure 19: DPMM Clustering Results for Gaussian data

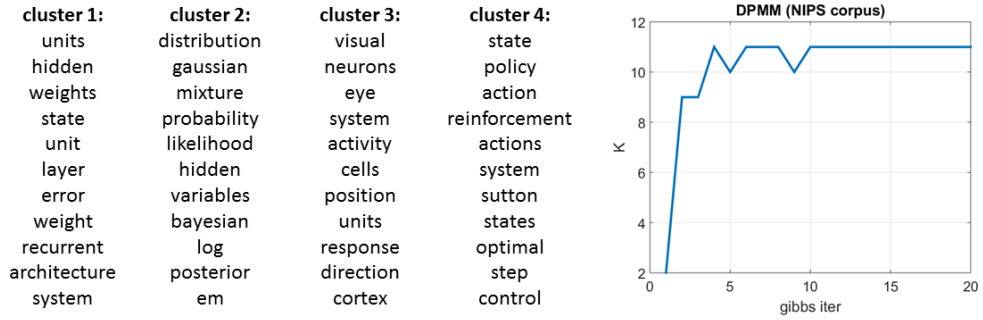


Figure 20: DPMM Clustering Results for Categorical data

and then make this approximation as close as possible to the true posterior $p^*(x)$. This reduces inference to an optimization problem.

We can use KL divergence to measure the distance between distributions. In particular, we use reverse KL to make the computation tractable.

$$KL(q||p^*) = \sum_x q(x) \log \frac{q(x)}{p^*(x)} \quad (64)$$

Let $\tilde{p}(x) = p^*(x)Z$ be the un-normalized distribution, then our objective function:

$$J(q) = KL(q||\tilde{p}) = \sum_x q(x) \log \frac{q(x)}{p^*(x)Z} = \sum_x q(x) \log \frac{q(x)}{p^*(x)} - \log Z = KL(q||p^*) - \log Z \quad (65)$$

Since KL divergence is non-negative, $J(q)$ is an upper bound on the marginal likelihood:

$$J(q) = KL(q||p^*) - \log Z \geq -\log Z = -\log p(D) \quad (66)$$

when $q(x)$ equals the true posterior $p^*(x)$, the KL divergence vanishes and the optimal value $J(q^*)$ equals the log partition function and for all other values of q it yields a bound. $J(q)$ is called the *variational free energy* and can be written as:

$$\min_q J(q) = E_q[\log q(x)] + E_q[-\log \tilde{p}(x)] = -H(q) + E_q[E(x)] \quad (67)$$

The variational objective function (67) is closely related to energy minimization in statistical physics. The first term acts as a regularizer by encouraging maximum entropy, while the second term is the expected energy and encourages the variational distribution q to explain the data.

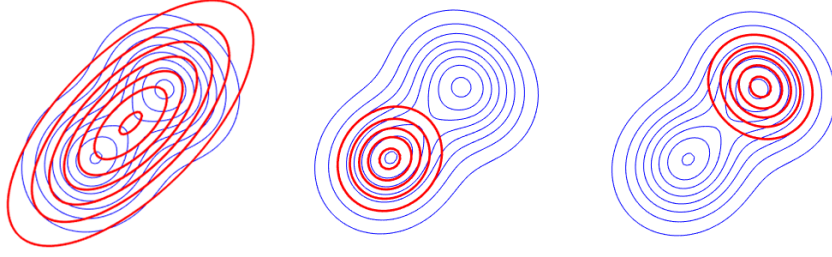


Figure 21: Forward vs reverse KL on a bimodal distribution. The blue contours is the true distribution $p(x)$. The red contours is approximate distribution $q(x)$.

The reverse KL that acts as a penalty term in the variational objective is also known as I-projection or information projection. In the reverse KL, $q(x)$ will typically under-estimate the support of $p(x)$ and will lock onto one of its modes. This is due to $q(x) = 0$ whenever $p(x) = 0$ to make sure the KL divergence stays finite. On the other hand, the forward KL, known as M-projection or moment projection is zero avoiding for $q(x)$ and will over-estimate the support of $p(x)$ as shown in Figure 21.

We can use the Jensen's inequality to derive the Evidence Lower Bound (ELBO), an objective that we can maximize to learn the variational parameters of our model. Let x be our data and z be the latent variables, then we can derive our ELBO objective as follows:

$$\log p(x) = \log \sum_z p(x, z) = \log \sum_z \frac{q(z)}{q(z)} p(x, z) = \log E_{q(z)} \left[\frac{p(x, z)}{q(z)} \right] \quad (68)$$

$$\geq E_{q(z)} \left[\log \frac{p(x, z)}{q(z)} \right] = E_{q(z)} \left[\log p(x, z) \right] - E_{q(z)} \left[\log q(z) \right] = \text{ELBO} \quad (69)$$

Notice that the first term is the average negative energy and the second term is the entropy. Thus, a good posterior must assign most of its probability mass to regions of low energy (i.e. high joint probability density) while also maximizing the entropy of $q(z)$. Thus, variational Bayes, in contrast to MAP, prevents $q(z)$ from collapsing to an atom.

One form of ELBO emphasizes that the lower bound becomes tighter as the variational distribution better approximates the posterior:

$$\text{ELBO} = E_{q(z)} \left[\log \frac{p(x, z)}{q(z)} \right] = E_{q(z)} \left[\log \frac{p(z|x)p(x)}{q(z)} \right] = -KL(q(z)||p(z|x)) + \log p(x) \quad (70)$$

Therefore, we can improve the ELBO by improving the model log evidence $\log p(x)$ through the prior $p(z)$ or the likelihood $p(z|x)$ or by improving the variational posterior approximation $q(z)$.

Finally, we can write the ELBO as follows:

$$\text{ELBO} = \frac{1}{n} \sum_{i=1}^n \left[E_{q(z_i)} [\log p(x_i|z_i)] - KL(q(z_i)||p(z_i)) \right] \quad (71)$$

this version emphasizes a likelihood term for the i -th observation and KL divergence term between each approximating distribution and the prior.

One of the most popular forms of variational inference is called the *mean field* approximation, where we assume that the posterior is a fully factorized approximation of the form:

$$q(x) = \prod_i q_i(x_i) \quad (72)$$

where we optimize over the parameters of each marginal distribution $q_j(x_j)$. Our goal is to minimize variational free energy $J(q)$ or equivalently, maximize the lower bound:

$$L(q) = -J(q) = \sum_x q(x) \log \frac{\tilde{p}(x)}{q(x)} \quad (73)$$

We can re-write the objective for each marginal distribution q_j , keeping the rest of the terms as constants:

$$L(q_j) = \sum_x \prod_i q_i(x_i) \left[\log \tilde{p}(x) - \sum_k \log q_k(x_k) \right] \quad (74)$$

$$= \sum_{x_j} \sum_{x_{-j}} q_j(x_j) \prod_{i \neq j} q_i(x_i) \left[\log \tilde{p}(x) - \sum_k \log q_k(x_k) \right] \quad (75)$$

$$= \sum_{x_j} q_j(x_j) \log f_j(x_j) - \sum_{x_j} q_j(x_j) \sum_{x_{-j}} \prod_{i \neq j} q_i(x_i) \left[\sum_{k \neq j} \log q_k(x_k) + \log q_j(x_j) \right] \quad (76)$$

$$= \sum_{x_j} q_j(x_j) \log f_j(x_j) - \sum_{x_j} q_j(x_j) \log q_j(x_j) + \text{const} \quad (77)$$

where we defined $\log f_j(x_j) = \sum_{x_{-j}} \prod_{i \neq j} q_i(x_i) \log \tilde{p}(x) = E_{-q_j}[\log \tilde{p}(x)]$. Since we are replacing the values by their mean value, the method is known as mean field. We can re-write $L(q_j) = -KL(q_j || f_j)$ and therefore maximize the objective by setting $q_j = f_j$ or equivalently:

$$\log q_j(x_j) = \log f_j(x_j) = E_{-q_j}[\log \tilde{p}(x)] \quad (78)$$

where the functional form of q_j will be determined by the type of variables x_j and their probability model.

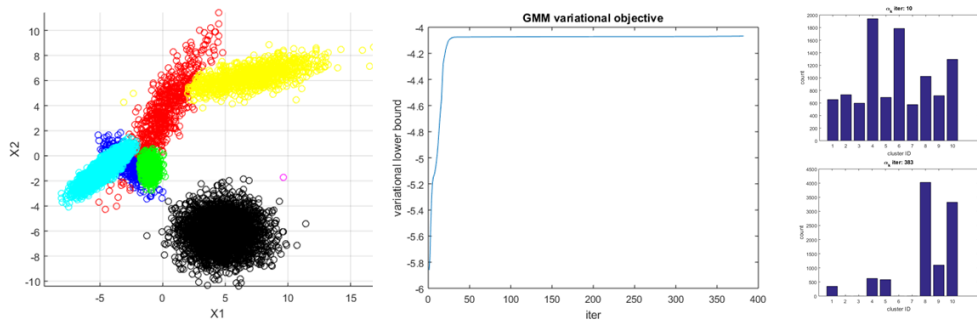


Figure 22: Variational Bayes EM applied to a Mixture of Gaussians

Figure 22 shows the posterior clustering assignment as a result of running Variational Bayes EM algorithm on a Gaussian Mixture. The algorithm correctly identified $K = 6$ clusters using 10 clusters as the starting point. We can see the rapid increase in the variational objective when the algorithm figures out that it could increase the objective by removing unnecessary clusters in early iterations, while the plateau results in moving the clusters around. This is also evident from the prior and posterior parameter values for mixing proportions α_k . As the number of iterations increase the counts for unnecessary mixture components drop to zero.

2.1 Application: Topic Models

A topic model is a latent variable model for discrete data such as text documents. Latent Dirichlet Allocation (LDA) is a topic model that represents each document as a finite mixture of topics, where a topic is a distribution over words. The objective is to learn the shared topic distribution and topic proportions for each document. LDA assumes a bag of words model in which the words are exchangeable and as a result sentence structure is not preserved, i.e. only the word counts matter. Thus, each document is reduced to a vector of counts over the vocabulary V and the entire corpus of D documents is summarized in a *term-document* matrix $A_{V \times D}$. LDA can be seen as a non-negative matrix factorization problem that takes the term-document matrix and factorizes it into a product of topics $W_{V \times K}$ and topic proportions $H_{K \times D}$: $A = WH$.

A common method for adjusting the word counts is *tf-idf* that logarithmically drives down to zero word counts that occur frequently across documents: $A_{t,d} \log \frac{D}{n_t}$, where D is the total number of documents in the corpus and n_t is the number of documents where term t appears. The *tf-idf* smoothing identifies the sets of words that are discriminative for documents and leads to better model performance. The term-document matrix generalizes from counts of individual words (unigrams) to larger structural units such as n -grams. In the case of n -grams different smoothing techniques (such as Laplace smoothing) are used to address the lack of observations in a very large feature space.

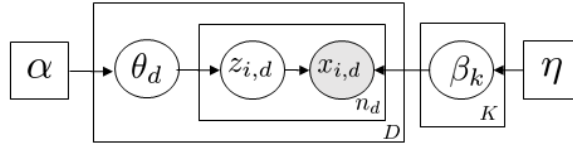


Figure 23: Latent Dirichlet Allocation (LDA) graphical model

Figure 23 shows the graphical model for the Latent Dirichlet Allocation (LDA). The LDA topic model associates each word $x_{i,d}$ with a topic label $z_{i,d} \in \{1, 2, \dots, K\}$. Each document is associated with topic proportions θ_d that could be used to measure document similarity. The topics β_k are shared across all documents. The hyper-parameters α and η capture our prior knowledge of topic proportions and topics, respectively, e.g. from past on-line training of the model. The full generative model can be specified as follows:

$$\theta_d | \alpha \sim \text{Dir}(\alpha) \quad (79)$$

$$z_{i,d} | \theta_d \sim \text{Cat}(\theta_d) \quad (80)$$

$$\beta_k | \eta \sim \text{Dir}(\eta) \quad (81)$$

$$x_{i,d} | z_{i,d} = k, \beta \sim \text{Cat}(\beta_k) \quad (82)$$

The joint distribution for a single document d can be written as follows [4]:

$$p(x, z, \theta, | \alpha, \beta) = p(\theta_d | \alpha) \prod_{i=1}^{n_d} p(z_{i,d} | \theta_d) p(x_{i,d} | z_{i,d}, \beta) \quad (83)$$

The parameters α and β are corpus-level parameters, the variable θ_d is sampled once every document, while $z_{i,d}$ and $x_{i,d}$ are word-level variables sampled once for each word in each document. Unlike a multinomial clustering model where each document is associated with a single topic,

LDA represents each document as a mixture of topics.

The key inference problem that we need to solve in order to use LDA is that of computing the posterior distribution of the latent variables for a given document: $p(\theta, z|x, \alpha, \beta)$. The posterior can be approximated with the following variational distribution:

$$q(\theta, z|\gamma, \phi) = q(\theta|\gamma) \prod_{i=1}^n q(z_i|\phi_i) \quad (84)$$

The variational parameters are optimized to maximize the Evidence Lower Bound (ELBO):

$$\log p(x|\alpha, \eta) \geq L(x, \phi, \gamma, \lambda) = E_q[\log p(x, z, \theta, \beta|\alpha, \eta)] - E_q[\log q(z, \theta, \beta)] \quad (85)$$

We choose a fully factored distribution q of the form:

$$q(z_{id} = k) = \phi_{dwk}; \quad q(\theta_d) \sim \text{Dir}(\theta_d|\gamma_d); \quad q(\beta_k) \sim \text{Dir}(\beta_k|\lambda_k) \quad (86)$$

We can expand the lower bound by using the factorizations of p and q :

$$L(\gamma, \phi; \alpha, \beta) = \mathbb{E}_q[\log p(\theta|\alpha)] + \mathbb{E}_q[\log p(z|\theta)] + \mathbb{E}_q[\log p(w|z, \beta)] - \mathbb{E}_q[\log q(\theta)] - \mathbb{E}_q[\log q(z)]$$

Each of the five terms in $L(\gamma, \phi; \alpha, \beta)$ can be expanded [4] as follows:

$$L(\gamma, \phi; \alpha, \beta) = \log \Gamma(\sum_{j=1}^k \alpha_j) - \sum_{i=1}^k \log \Gamma(\alpha_i) + \sum_{i=1}^k (\alpha_i - 1)(\Psi(\gamma_i) - \Psi(\sum_{j=1}^k \gamma_j)) \quad (87)$$

$$+ \sum_{n=1}^N \sum_{i=1}^k \phi_{ni} (\Psi(\gamma_i) - \Psi(\sum_{j=1}^k \gamma_j)) \quad (88)$$

$$+ \sum_{n=1}^N \sum_{i=1}^k \sum_{j=1}^V \phi_{ni} w_n^j \log \beta_{ij} \quad (89)$$

$$- \log \Gamma(\sum_{j=1}^k \gamma_j) + \sum_{i=1}^k \log \Gamma(\gamma_i) - \sum_{i=1}^k (\gamma_i - 1)(\Psi(\gamma_i) - \Psi(\sum_{j=1}^k \gamma_j)) \quad (90)$$

$$- \sum_{n=1}^N \sum_{i=1}^k \phi_{ni} \log \phi_{ni} \quad (91)$$

where $\Psi(x) = \frac{d}{dx} \log \Gamma(x)$ is the digamma function. $L(\gamma, \phi; \alpha, \beta)$ can be maximized using coordinate ascent over the variational parameters ϕ, γ, λ [4]:

$$\phi_{dwk} \propto \exp\{E_q[\log \theta_{dk}] + E_q[\log \beta_{kw}]\} \quad (92)$$

$$\gamma_{dk} = \alpha + \sum_w n_{dw} \phi_{dwk} \quad (93)$$

$$\lambda_{kw} = \eta + \sum_d n_{dw} \phi_{dwk} \quad (94)$$

where the expectations under q of $\log \theta$ and $\log \beta$ are:

$$E_q[\log \theta_{dk}] = \Psi(\gamma_{dk}) - \Psi(\sum_{i=1}^K \gamma_{di}) \quad E_q[\log \beta_{kw}] = \Psi(\lambda_{kw}) - \Psi(\sum_{i=1}^W \lambda_{ki}) \quad (95)$$

$W_{V \times K}$ where only the first 64 rows are shown. The perplexity plots show the improvement in model ability to explain the data as the number of EM iterations increases, where perplexity is defined as follows:

$$\text{Perplexity}(w_{\text{test}}) = \exp\left\{-\frac{1}{D_{\text{test}}} \sum_d \frac{1}{n_d} \sum_{w \in n_d} \log p(w_{\text{test}})\right\} \quad (96)$$

Model selection can be done by evaluating perplexity for different values of K . Guided by the fact that there are 20 different newsgroups, the value of K was set to 20. Finally, the top 20 words for a random sample of 4 topics are shown in Figure 24. The four topics are about information, sports, computers and school.

2.2 Application: Image Denoising in Ising Model

The Ising model is an example of a Markov Random Field (MRF) and has its origins in statistical physics. The Ising model assumes that we have a grid of nodes, where each node can be in one of two possible states. The state of each node depends on the neighboring nodes through interaction potentials. In the case of images, this translates to a smoothness constraint, i.e. a pixel prefers to be of the same color as the neighboring pixels. In the image denoising problem, we assume that we have a 2-D grid of noisy pixel observations of an underlying true image and we would like to recover the true image.

Let y_i be noisy observations of binary latent variables $x_i \in \{-1, +1\}$. We can write down the joint distribution as follows:

$$p(x, y) = p(x)p(y|x) = \prod_{(s,t) \in E} \Psi_{st}(x_s, x_t) \prod_{i=1}^n p(y_i|x_i) = \prod_{(s,t) \in E} \exp\{x_s w_{st} x_t\} \prod_{i=1}^n N(y_i|x_i, \sigma^2) \quad (97)$$

where the interaction potentials are represented by Ψ_{st} for every pair of nodes x_s and x_t in a set of edges E and the observations y_i are Gaussian with mean x_i and variance σ^2 . Here, w_{st} is the coupling strength and assumed to be constant and equal to $J > 0$ indicating a preference for the same state as neighbors (i.e. the potential $\Psi(x_s, x_t) = \exp\{x_s J x_t\}$ is higher when x_s and x_t are both either $+1$ or -1).

To fit the model parameters using variational inference, we want to maximize the ELBO:

$$\text{ELBO} = E_{q(x)} \left[\log p(x, y) \right] - E_{q(x)} \left[\log q(x) \right] = \quad (98)$$

$$= E_{q(x)} \left[\sum_{(s,t) \in E} x_s w_{st} x_t + \sum_{i=1}^n \log N(x_i; \sigma^2) \right] - \sum_{i=1}^n E_{q_i(x)} \left[\log q_i(x) \right] \quad (99)$$

where we are using the mean-field assumption of a fully-factored approximation $q(x)$:

$$q(x) = \prod_{i=1}^n q(x_i; \mu_i) \quad (100)$$

Using the previously derived result, we state that $q(x_i; \mu_i)$ that minimizes the KL divergence is given by:

$$q_i(x_i) = \frac{1}{Z_i} \exp \left[E_{-q_i} \{ \log p(x) \} \right] \quad (101)$$

where E_{-q_i} denotes the expectation over every q_j except for $j = i$. To compute $q_i(x_i)$, we only care about the terms that involve x_i , i.e. we can isolate them as follows:

$$E_{-q_i}\{\log p(x)\} = E_{-q_i}\{x_i \sum_{j \in N(i)} w_{ij}x_j + \log N(x_i, \sigma^2) + \text{const}\} = \quad (102)$$

$$= x_i \sum_{j \in N(i)} J \times \mu_j + \log N(x_i, \sigma^2) + \text{const} \quad (103)$$

where $N(i)$ denotes the neighbors of node i and μ_j is the mean of a binary random variable:

$$\mu_j = E_{q_j}[x_j] = q_j(x_j = +1) \times (+1) + q_j(x_j = -1) \times (-1) \quad (104)$$

Figure 25 shows the parametric form of our mean-field approximation for the Ising model:

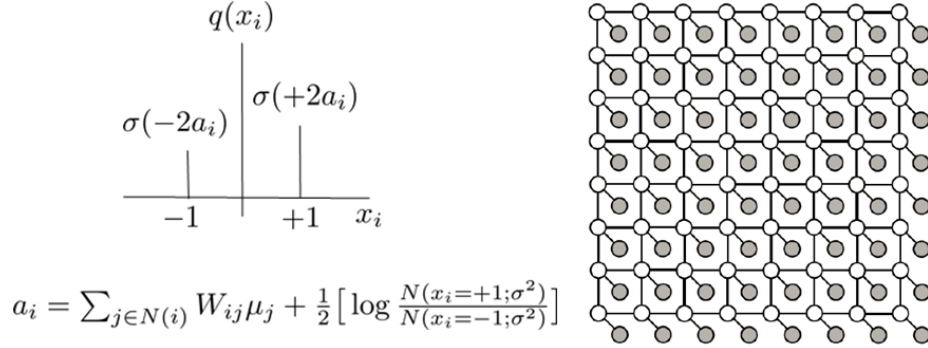


Figure 25: Mean-field VI approximation (left) and Ising model (right).

In order to compute this mean, we need to know the values of $q_j(x_j = +1)$ and $q_j(x_j = -1)$. Let $m_i = \sum_{j \in N(i)} w_{ij} \mu_j$ be the mean value of neighbors and let $L_i^+ = N(x_i = +1, \sigma^2)$ and $L_i^- = N(x_i = -1, \sigma^2)$, then we can compute the mean as follows:

$$q_i(x_i = +1) = \frac{\exp\{m_i + L_i^+\}}{\exp\{m_i + L_i^+\} + \exp\{-m_i + L_i^-\}} \quad (105)$$

$$= \frac{1}{1 + \exp\{-2m_i + L_i^- - L_i^+\}} = \frac{1}{1 + \exp\{-2a_i\}} = \sigma(2a_i) \quad (106)$$

where $a_i = m_i + 1/2(L_i^+ - L_i^-)$ and $\sigma(x)$ is a sigmoid function. Since $q_i(x_i = -1) = 1 - q_i(x_i = +1) = 1 - \sigma(2a_i) = \sigma(-2a_i)$, we can write the mean of our variational approximation $q_i(x_i)$ as follows:

$$\mu_i = E_{q_i}[x_i] = \sigma(2a_i) - \sigma(-2a_i) = \tanh(a_i) \quad (107)$$

In other words, our mean-field updates of the variational parameters μ_i at iteration k are computed as follows:

$$\mu_i^{(k)} = \tanh \left(\sum_{j \in N(i)} w_{ij} \mu_j^{(k-1)} + \frac{1}{2} \left[\log \frac{N(x_i = +1, \sigma^2)}{N(x_i = -1, \sigma^2)} \right] \right) \times \lambda + (1 - \lambda) \times \mu_i^{(k-1)} \quad (108)$$

where we added a learning rate parameter $\lambda \in (0, 1]$. We further note that we can simplify the

computation of ELBO term by term as follows:

$$\sum_{(s,t) \in E} E_{q(x)}[x_s w_{st} x_t] = \frac{1}{2} \sum_{i=1}^n \sum_{j \in N(i)} \left(\sum_{x_i \in \{-1, +1\}} \sum_{x_j \in \{-1, +1\}} q_i(x_i) q_j(x_j) x_i J x_j \right) = \quad (109)$$

$$\frac{1}{2} \sum_{i=1}^n \sum_{j \in N(i)} \left(q_i(x_i = +1) J E[x_j] - q_i(x_i = -1) J E[x_j] \right) = \frac{1}{2} \sum_{i=1}^n \sum_{j \in N(i)} E[x_i] J E[x_j] \quad (110)$$

Similarly,

$$E_{q(x)}[\log N(x_i, \sigma^2)] = \sum_{i=1}^n \left[\sum_{x_i \in \{-1, +1\}} q_i(x_i) \log N(x_i, \sigma^2) \right] = \quad (111)$$

$$\sum_{i=1}^n \left[\sigma(2a_i) \log N(x_i = +1, \sigma^2) + \sigma(-2a_i) \log N(x_i = -1, \sigma^2) \right] \quad (112)$$

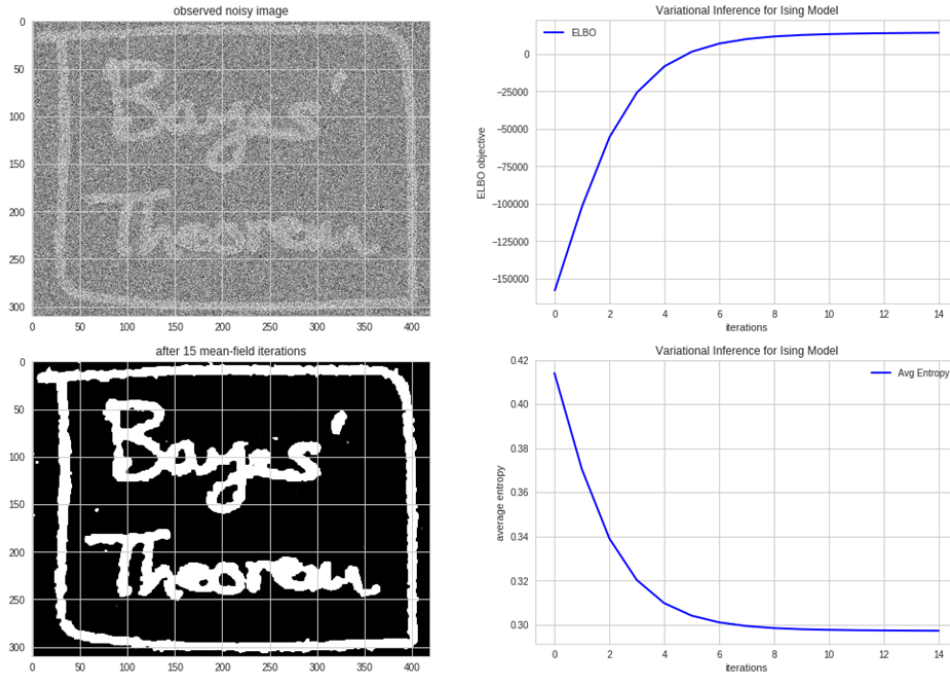


Figure 26: Image Denoising for Ising model.

Figure 26 shows experimental results for binary image denoising via mean-field variational inference. The noisy observed image is shown in the top-left obtained by adding Gaussian noise to each pixel and binarizing the image based on mean threshold. We then set the variational inference parameters such as the coupling strength $J = 1$, noise level $\sigma = 2$, smoothing rate $\lambda = 0.5$ and max number of iterations to 15. The resulting de-noised binary image is shown in the bottom-left corner of Figure 26. We can also see an increase in the ELBO objective (top-right) and a decrease in the average entropy of our binary random variables $q_i(x_i)$ representing the value of each pixel (bottom-right) as the number of mean-field iterations increases.

The 2-D Ising model can be extended in multiple ways, for example: 3-D grids and K -states per node (aka Potts model).

2.3 Stochastic Variational Inference

One limitation of LDA is that the number of topics is fixed ahead of time. A commonly used approach to finding the number of topics K is cross-validation. However, for very large data-sets this approach may not be practical. We can address this issue with a Bayesian non-parametric topic model where the number of topics is learned from data: the Hierarchical Dirichlet Process (HDP) topic model.

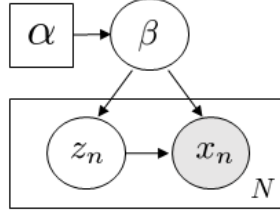


Figure 27: Graphical model with local and global latent variables.

While traditional algorithms require repeatedly analyzing the whole data set before updating the variational parameters, we focus on analyzing randomly sampled subsets. First, we derive the Stochastic Variational Inference (SVI) algorithm for a class of graphical models with global and local latent variables as shown in Figure 27. The joint distribution factorizes into a product of a global term β and local term z_n :

$$p(x, z, \beta | \alpha) = p(\beta | \alpha) \prod_{n=1}^N p(x_n, z_n | \beta) \quad (113)$$

Our goal is to approximate the posterior distribution of the hidden variables given the observations $p(\beta, z | x)$. We use a variational distribution to approximate the posterior as measured by KL divergence. Variational inference minimizes KL divergence or alternatively maximizes the evidence lower bound (ELBO):

$$\log p(x) = \log \int p(x, z, \beta) dz d\beta \quad (114)$$

$$= \log \int p(x, z, \beta) \frac{q(z, \beta)}{q(z, \beta)} dz d\beta \quad (115)$$

$$= \log \left(E_q \left[\frac{p(x, z, \beta)}{q(z, \beta)} \right] \right) \quad (116)$$

$$\geq E_q[\log p(x, z, \beta)] - E_q[\log q(z, \beta)] \quad (117)$$

$$= L(q) \quad (118)$$

The ELBO contains two terms: the first term is the expected log joint $E_q[\log p(x, z, \beta)]$ and the second term is the entropy of the variational distribution $-E_q[\log q(z, \beta)]$. We restrict $q(z, \beta)$ to be in a tractable family of distributions in order to efficiently compute the expectations in the ELBO. We then find the member of the family that maximizes the ELBO and use the optimized distribution as a proxy for the posterior. Solving this maximization problem is equivalent to finding the member of the family that is closest in KL divergence to the posterior:

$$\min KL(q(z, \beta) || p(z, \beta | x)) = E_q[\log q(z, \beta)] - E_q[\log p(z, \beta | x)] \quad (119)$$

$$= E_q[\log q(z, \beta)] - E_q[\log p(x, z, \beta)] + \log p(x) \quad (120)$$

$$= -L(q) + \text{const} \quad (121)$$

where $\log p(x)$ is replaced by a constant because it does not depend on q . The simplest variational family of distributions is the fully factored mean-field family. In this family, each hidden variable is independent and governed by its own parameter:

$$q(z, \beta) = q(\beta|\lambda) \prod_{n=1}^N \prod_{j=1}^J q(z_{nj}|\phi_{nj}) \quad (122)$$

To specify the form of the distribution, we choose $q(\beta|\lambda)$ and $q(z_{nj}|\phi_{nj})$ to be in the exponential distribution with the natural parameters λ and ϕ_{nj} :

$$q(\beta|\lambda) = h(\beta) \exp\{\lambda^T t(\beta) - a_g(\lambda)\} \quad (123)$$

$$q(z_{nj}|\phi_{nj}) = h(z_{nj}) \exp\{\phi_{nj}^T t(z_{nj}) - a_l(\phi_{nj})\} \quad (124)$$

The mean-field family has several computational advantages. For example the entropy term in the ELBO objective function decomposes:

$$-E_q[\log q(z, \beta)] = -E_\lambda[\log q(\beta)] - \sum_{n=1}^N \sum_{j=1}^J E_{\phi_{nj}}[\log q(z_{nj})] \quad (125)$$

In traditional mean-field variational inference, we maximize $L(q)$ with coordinate ascent: we iteratively optimize each variational parameter while holding the other parameters fixed. Given our assumptions that each conditional is an exponential family, we can optimize each coordinate in closed form. We first derive the coordinate update for the global parameter λ . Keeping only the terms of $L(q)$ that depend on λ , we get:

$$L(\lambda) = E_q[\log p(x, z, \beta)] - E_q[\log q(\beta)] + \text{const} \quad (126)$$

$$= E_q[\log p(\beta|x, z)] - E_q[\log q(\beta)] + \text{const} \quad (127)$$

where we used the factorization $p(x, z, \beta) = p(\beta|x, z)p(x, z)$ and absorbed $E_q[p(x, z)]$ into the constant that does not depend on λ . To derive the coordinate ascent update, we take the gradient [19]:

$$\nabla_\lambda L = \nabla_\lambda^2 a_g(\lambda)(E_q[\eta_g(x, z, \alpha)] - \lambda) \quad (128)$$

We can set the gradient to zero by setting $\lambda = E_q[\eta_g(x, z, \alpha)]$. This sets the global variational parameter equal to the expected natural parameter of its complete conditional distribution. We now turn to the local parameters ϕ_{nj} . The gradient is similar to the global case:

$$\nabla_{\phi_{nj}} L = \nabla_{\phi_{nj}}^2 a_l(\phi_{nj})(E_q[\eta_l(x_n, z_{n,-j}, \beta)] - \phi_{nj}) \quad (129)$$

We can set the gradient to zero by choosing $\phi_{nj} = E_q[\eta_l(x_n, z_{n,-j}, \beta)]$. The variational updates form the algorithm for coordinate ascent variational inference, iterating between updates of each local parameter and the global parameter. The full algorithm is described in Algorithm 3 which is guaranteed to find a local optimum of the ELBO.

We now turn the Hierarchical Dirichlet Process (HDP) topic model. The HDP topic model couples a set of document-level DPs via a single top-level DP [39]. The base distribution H of the top-level DP is a symmetric Dirichlet over the vocabulary simplex. We draw once from this DP: $G_0 \sim DP(\omega, H)$. In the second level, we use G_0 as a base measure for a document level DP: $G_d \sim DP(\alpha, G_0)$. As a result, the global topics are shared between documents with different mixing proportions.

Figure 28 shows a graphical model for the HDP topic model. The generative process of the HDP topic model can be described as follows.

Algorithm 3 Coordinate Ascent SVI [19]

- 1: Initialize $\lambda^{(0)}$ randomly
 - 2: **repeat**
 - 3: **for** each local variational parameter ϕ_{nj} **do**
 - 4: Update $\phi_{nj}^{(t)} = E_{q^{(t-1)}}[\eta_{l,j}(x_n, z_{n,-j}, \beta)]$
 - 5: **end for**
 - 6: Update the global variational parameters: $\lambda^{(t)} = E_{q^{(t)}}[\eta_g(z, x)]$
 - 7: **until** the ELBO converges
-

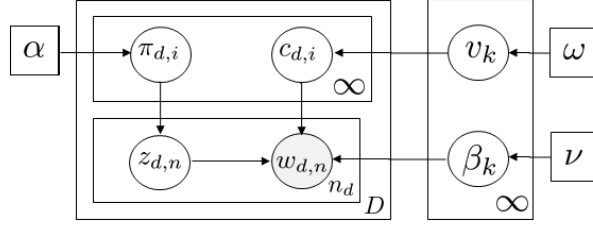


Figure 28: HDP graphical model for SVI inference.

1. Draw global topics, $\beta_k \sim \text{Dir}(\eta)$
2. Draw mixing proportions, $\nu_k \sim \text{Beta}(1, \omega)$
3. For each document d :
 - (a) Draw document-level topic indices, $c_{di} \sim \text{Mult}(\nu)$
 - (b) Draw document-level mixing proportions, $\pi_{di} \sim \text{Beta}(1, \alpha)$
 - (c) For each word n :
 - i. Draw topic assignment $z_{dn} \sim \text{Mult}(\pi_d)$
 - ii. Draw word $w_n \sim \text{Mult}(\beta_{c_d, z_{dn}})$

In order to implement an infinite number of topics at corpus and document levels we need to truncate our representation to K topics at the corpus level and T topics at the document level. This way we are optimizing a finite number of variational parameters. We can write down the joint variational family distribution as:

$$q(\beta, \nu, z, \pi) = \left(\prod_{k=1}^K q(\beta_k | \lambda_k) q(\nu_k | a_k) \right) \left(\prod_{d=1}^D \prod_{i=1}^T q(c_{di} | \zeta_{di}) q(\pi_{di} | \gamma_{di}) \prod_{n=1}^N q(z_{dn} | \phi_{dn}) \right) \quad (130)$$

The Stochastic Variational Inference (SVI) algorithm is summarized in Algorithm 4.

where the expectations used in Algorithm 4 are computed as follows:

$$E[\log \beta_{kv}] = \Psi(\lambda_{kv}) - \Psi\left(\sum_{v'} \lambda_{kv'}\right) \quad (131)$$

$$E[\log \sigma_k(V)] = \Psi(a_k) - \Psi(a_k + b_k) + \sum_{l=1}^{k-1} [\Psi(b_l) - \Psi(a_l + b_l)] \quad (132)$$

The on-line variational inference algorithm for Hierarchical Dirichlet Process was used to fit a topic model on the 20newsgroups dataset. The dataset consists of 11,314 documents and over

Algorithm 4 Stochastic Variational Inference for HDP [19]

- 1: Initialize $\lambda^{(0)}$ randomly. Set $a^{(0)} = 1$ and $b^{(0)} = \omega$
 - 2: Set the step-size schedule ρ_t
 - 3: **repeat**
 - 4: Sample a document w_d uniformly from the dataset
 - 5: For $i \in \{1, \dots, T\}$, $k \in \{1, \dots, K\}$: $\zeta_{di}^k \propto \exp\{\sum_{n=1}^N E[\log \beta_{k,w_{dn}}]\}$
 - 6: For $n \in \{1, \dots, N\}$, $i \in \{1, \dots, T\}$: $\phi_{dn}^i \propto \exp\{\sum_{k=1}^K \zeta_{di}^k E[\log \beta_{k,w_{dn}}]\}$
 - 7: **repeat**
 - 8: For $i \in \{1, \dots, T\}$ set

$$\begin{aligned} \gamma_{di}^{(1)} &= 1 + \sum_{n=1}^N \phi_{dn}^i \\ \gamma_{di}^{(2)} &= \alpha + \sum_{n=1}^N \sum_{j=i+1}^T \phi_{dn}^j \\ \zeta_{di}^k &\propto \exp\{E[\log \sigma_k(V)] + \sum_{n=1}^N \phi_{dn}^i E[\log \beta_{k,w_{dn}}]\}, k \in \{1, \dots, K\} \end{aligned}$$
 - 9: For $n \in \{1, \dots, N\}$ set

$$\phi_{dn}^i \propto \exp\{E[\log \sigma_i(\pi_d)] + \sum_{k=1}^K \zeta_{di}^k E[\log \beta_{k,w_{dn}}]\}, i \in \{1, \dots, T\}$$
 - 10: **until** local parameters converge
 - 11: For $k \in \{1, \dots, K\}$ set intermediate topics:

$$\begin{aligned} \hat{\lambda}_{kv} &= \eta + D \sum_{i=1}^T \zeta_{di}^k \sum_{n=1}^N \phi_{dn}^i w_{dn} \\ \hat{a}_k &= 1 + D \sum_{i=1}^T \zeta_{di}^k \\ \hat{b}_k &= \omega + D \sum_{i=1}^T \sum_{l=k+1}^K \zeta_{di}^l \end{aligned}$$
 - 12: Set

$$\begin{aligned} \lambda^{(t)} &= (1 - \rho_t) \lambda^{(t-1)} + \rho_t \hat{\lambda} \\ a^{(t)} &= (1 - \rho_t) a^{(t-1)} + \rho_t \hat{a} \\ b^{(t)} &= (1 - \rho_t) b^{(t-1)} + \rho_t \hat{b} \end{aligned}$$
 - 13: **until** end of documents
-

100K unique tokens. Standard text pre-processing was used including tokenization, stop-word removal and stemming. A compressed dictionary of 4K words was constructed by filtering out tokens that appear in less than 5 documents and more than 50% of the corpus. The top-level truncation was set to $T = 20$ topics and the second level truncation was set to $K = 8$ topics. The concentration parameters were chosen as $\gamma = 1.0$ at the top-level and $\alpha = 0.1$ at the group level to yield a broad range of shared topics that are concentrated at the group level. Figure 29 shows a sample of the global level topics inferred by online variational HDP algorithm. We can find topics about autos, politics and for sale items that correspond to the target labels of the 20newsgroups dataset.



Figure 29: Sample of HDP topics inferred using online variational bayes algorithm on 20news-groups dataset.

3 Optimization

3.1 Simulated Annealing

Simulated annealing is a stochastic algorithm that attempts to find the global optimum of an objective function $f(x)$. The method is inspired by statistical physics, in particular the Boltzmann distribution that specifies the probability of being in a particular state x :

$$p(x) \propto \exp\{-f(x)/T\} \quad (133)$$

where $f(x)$ is the energy of the system and T is the temperature. As the temperature approaches zero, the system spends more and more time in its minimum energy (most probable) state. As the temperature decreases, the largest peaks become larger and the smallest peaks disappear. By cooling slowly enough, it is possible to track the largest peak and therefore find the global optimum.

Simulated annealing is closely related to the Metropolis-Hastings algorithm for generating samples from a probability distribution. At each step of the algorithm, we sample a new state according to a proposal distribution $x' \sim q(\cdot|x_k)$, such as a random walk proposal:

$$x' = x_k + \epsilon_k, \quad \text{where } \epsilon_k \sim N(0, \Sigma) \quad (134)$$

Having proposed a new state, we compute α as in Algorithm 5. Thus, if a new state has lower energy (higher probability), we will definitely accept it but if it has higher energy (lower probability), we might still accept it depending on the temperature. Therefore, the algorithm allows downhill moves in probability space but less frequently as the temperature drops. In practice it is common to use an exponential cooling schedule: $T_k = T_0 C^k$, where $T_0 \sim 1$ is the initial temperature and $C \sim 0.8$ is the cooling rate. Cooling too quickly can result in getting stuck in local optima, while cooling too slowly wastes time. The optimum cooling schedule is difficult to determine.

Algorithm 5 Simulated Annealing

```
1:  $\alpha = \exp\{(f(x) - f(x'))/T\}$ 
2:  $r = \min(1, \alpha)$ 
3:  $u \sim \text{Unif}(0, 1)$ 
4: if  $u < r$ 
5:    $x_{k+1} = x'$ 
6: else
7:    $x_{k+1} = x_k$ 
8: end if
```

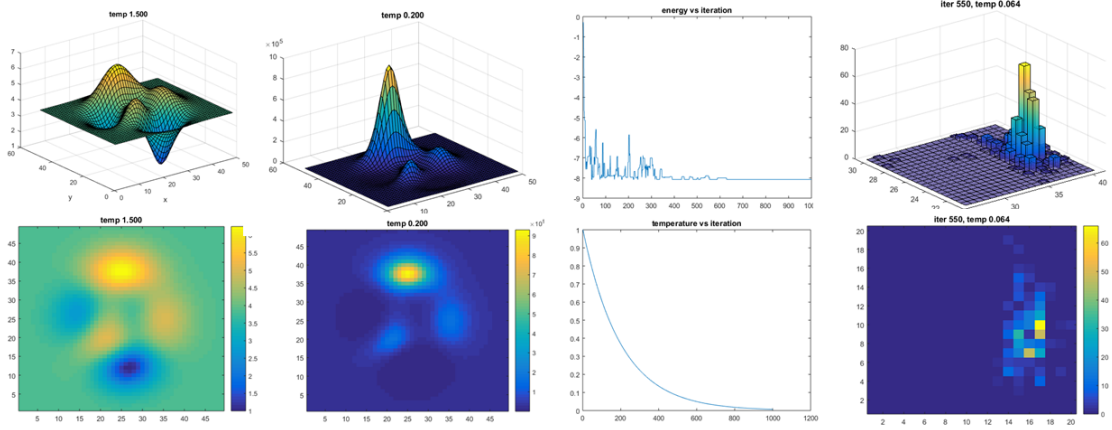


Figure 30: Simulated Annealing

Figure 30 shows the objective $f(x)$ at two different temperatures (left). The function appears more peaky when the temperature is lower. We can also see that the method stochastically reduces the energy over time for the given cooling schedule (middle). Finally a histogram of samples shows that most samples are concentrated near the global maximum (right).

3.2 Bayesian Optimization

Machine learning algorithms frequently require careful tuning of hyperparameters. Often exhaustive and computationally expensive methods such as grid search cross validation are used to find model parameters that optimize a suitable performance objective. An alternative to grid search is randomized parameter optimization that samples parameter settings from a distribution over possible parameter values. This has two main benefits over the exhaustive search: the number of iterations can be chosen independent of the number of parameters and adding parameters that do not influence performance does not decrease efficiency. Rather than exploring the parameter space randomly (according to a chosen distribution), it would be great to adapt an active learning approach that selects parameter values in a way that reduces uncertainty and provides a balance between exploration and exploitation. Bayesian optimization provides an automated Bayesian framework by utilizing Gaussian Processes (GPs) to model algorithm’s generalization performance [37].

Bayesian optimization assumes that a suitable performance function was sampled from a Gaussian Process and maintains a posterior distribution for this function as observations are made: $f(x) \sim \text{GP}(m(x), \kappa(x, x'))$. To choose which hyperparameters to explore next, one can optimize the Expected Improvement (EI) over the current best result or the Gaussian process

Upper Confidence Bound (UCB). EI and UCB have been shown to be efficient in the number of function evaluations required to find the global optimum of multi-modal black-box functions. Bayesian optimization uses all of the information available from previous evaluations of the objective function as opposed to relying on local gradient and Hessian approximations. This results in an automated procedure that can find an optimum of non-convex functions with relatively few evaluations, at the cost of performing more computation to determine the next point to try. This is particularly useful when evaluations are expensive to perform such as in selecting hyperparameters for deep neural networks.

To determine what point should be evaluated next, we need to choose an acquisition function which is used to construct a utility function from the GP posterior. In general, the acquisition function depends on previous observations as well as GP hyperparameters that we denote as $a(x; \{x_n, y_n\}, \theta)$, then $x_{next} = \arg \max_x a(x)$. Let $\mu(x; \{x_n, y_n\}, \theta)$ be the predictive GP mean function, $\sigma^2(x; \{x_n, y_n\}, \theta)$ be the predictive GP variance function and $\Phi(x)$ be the cumulative distribution function of the standard normal. Then we can define the following acquisition functions.

Probability of Improvement. This strategy maximizes the probability of improving over the best current value. This can be computed as follows:

$$a_{PI}(x; \{x_n, y_n\}, \theta) = \Phi(\gamma(x)), \quad \gamma(x) = \frac{f(x_{best}) - \mu(x; \{x_n, y_n\}, \theta)}{\sigma(x; \{x_n, y_n\}, \theta)} \quad (135)$$

Expected Improvement. Alternatively, one could choose to maximize the expected improvement (EI) over the current best.

$$a_{EI}(x; \{x_n, y_n\}, \theta) = \sigma(x; \{x_n, y_n\}, \theta)(\gamma(x)\Phi(\gamma(x)) + N(\gamma(x); 0, 1)) \quad (136)$$

Upper Confidence Bound. UCB is the idea of exploiting upper confidence bounds to construct acquisition functions that minimize regret over the course of their optimization. These acquisition functions have the following form:

$$a_{UCB}(x; \{x_n, y_n\}, \theta) = \mu(x; \{x_n, y_n\}, \theta) - \kappa \sigma(x; \{x_n, y_n\}, \theta) \quad (137)$$

where κ is a tunable acquisition parameter to balance exploration and exploitation. The optima of acquisition functions are located where the uncertainty in GP posterior is large (exploration) and/or where the GP prediction is high (exploitation). Since acquisition functions have analytical forms that are simple to evaluate, they are easier to optimize than the original objective function.

An important practical consideration for Bayesian optimization is an appropriate choice of GP kernel and its associated hyperparameters. Squared exponential kernel is often a default choice for Gaussian process regression. However, sample functions with this kernel are too smooth for practical optimization problems. Instead the following ARD Matern 5/2 kernel is proposed:

$$K_{M52}(x, x') = \theta_0 \left(1 + \sqrt{5r^2(x, x')} + \frac{5}{3}r^2(x, x') \right) \exp\{-\sqrt{5r^2(x, x')}\} \quad (138)$$

$$r^2(x, x') = \sum_{d=1}^D (x_d - x'_d)^2 / \theta_d^2 \quad (139)$$

This kernel results in sample functions that are twice-differentiable without requiring the smoothness of the squared exponential kernel. The associated kernel hyperparameters are D length scales $\theta_{1:D}$, the covariance amplitude θ_0 , the observation noise ν and a constant mean m . An additional assumption made by GP regression is that the underlying process is stationary, i.e. we can re-write the kernel $k(x, x')$ as a function of $x - x'$. Intuitively, a function whose length scale does not change throughout the input space will be modelled well by a GP with a stationary kernel.

Bayesian Optimization algorithm is summarized in Algorithm 6.

Algorithm 6 Bayesian Optimization

- 1: **for** $n = 1, 2, \dots$ **do**
 - 2: select new x_{n+1} by optimizing acquisition function α
 $x_{n+1} = \arg \max_x \alpha(x; D_n, \theta)$
 - 3: query objective function to obtain $y_{n+1} = f(x_{n+1})$
 - 4: augment data $D_{n+1} = \{D_n, (x_{n+1}, y_{n+1})\}$
 - 5: update GP posterior and acquisition function
 - 6: **end for**
-

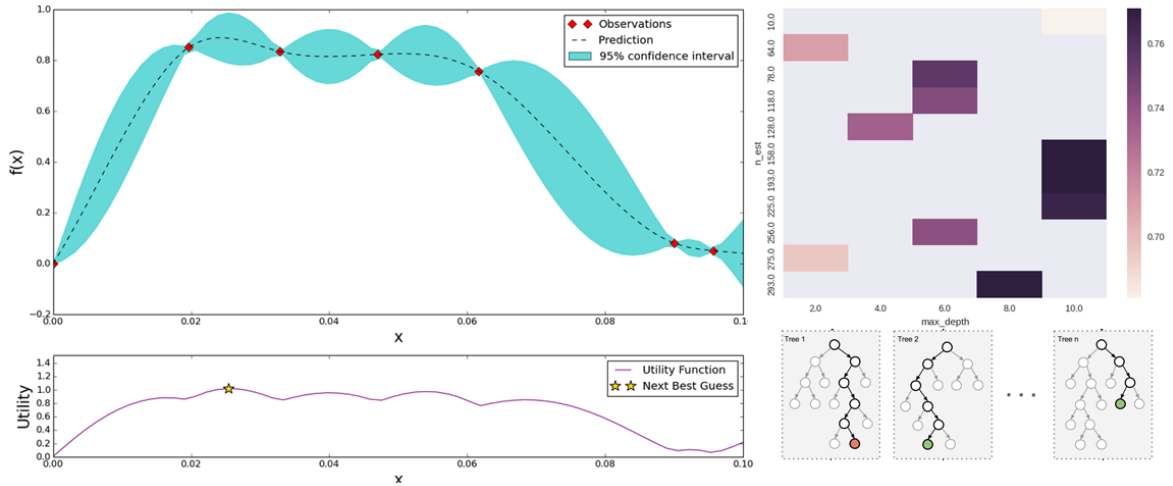


Figure 31: Bayesian Optimization of SVM and Random Forest parameters

Figure 31 shows Bayesian optimization applied to SVM and Random Forest. F1 score was used as performance objective function for a classification task. The figure on the left shows Bayesian optimization of F1 score as a function of the gamma parameter of the SVM RBF kernel: $K(x, x') = \exp\{-\gamma\|x - x'\|^2\}$. We can see that after only 7 iterations we have discovered the gamma parameter that gives the maximum F1 score. The peak of EI utility function at the bottom tells us which experiment to perform next. The figure on the right shows Bayesian optimization of F1 score as a function of maximum depth and the number of estimators of a Random Forest classifier. From the heatmap, we can tell that the maximum F1 score is achieved for 158 estimators with depth equal to 10.

3.3 Active Learning

The key idea behind active learning is that a machine learning algorithm can achieve greater accuracy with fewer training labels if it is allowed to choose the data from which it learns. Active learning is well motivated in many modern machine learning problems where unlabeled data may be abundant but labels are expensive to obtain. Active learning is sometimes called query learning or optimal experimental design because an active learner poses queries in the form of unlabelled data instances to be labeled by an oracle. In this way, the active learner seeks to achieve high accuracy using as few labeled instances as possible [35].

We focus on *pool-based sampling* that assumes that there is a small set of labeled data L and a large pool of unlabelled data U . Queries are selectively drawn from the pool according to an informativeness measure. Pool based methods rank the entire collection of unlabelled data to select the best query. Therefore, for very large data-sets, stream-based sampling maybe more appropriate where the data is scanned sequentially and query decisions are evaluated individually.

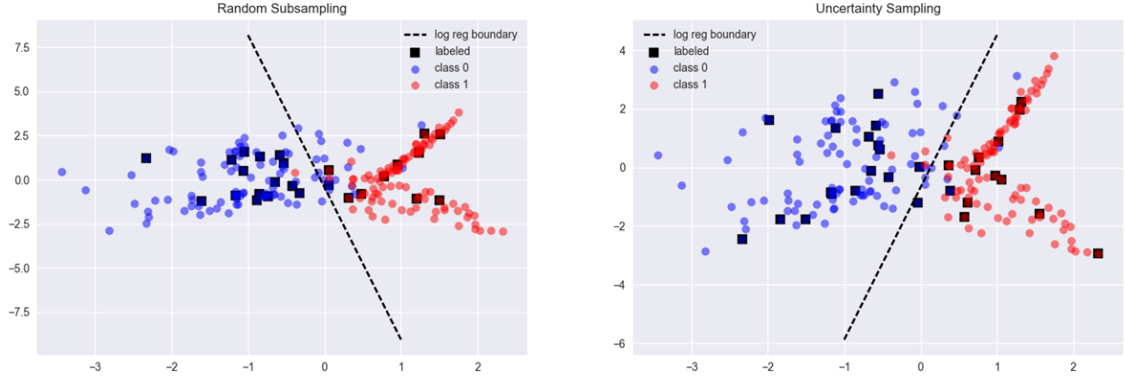


Figure 32: Random Subsampling vs Uncertainty Sampling

Figure 32 shows an example of pool-based active learning based on binary classification of a synthetic dataset with two balanced classes using Logistic Regression (LR). On the left, we can see the LR decision boundary as a result of training on a randomly subsampled set of 30 labels that achieves a classification accuracy of 90% on held out data. On the right, we can see the LR decision boundary as a result of training on 30 queries selected by uncertainty sampling based on entropy. Uncertainty sampling achieves a higher classification accuracy of 92.5% on the held out set.

3.3.1 Query Strategies

Uncertainty Sampling. One of the simplest and most commonly used query framework is uncertainty sampling. In this framework, an active learner queries the label about which it is least certain. For example, in binary logistic regression, uncertainty sampling queries points near the boundary where the probability of being positive is close to 0.5. For multi-class problems, uncertainty sampling can query points that are least confident:

$$x_{LC}^* = \arg \max_x 1 - P_{\theta}(\hat{y}|x) \quad (140)$$

where $\hat{y} \in \{1, \dots, K\}$ is the class label with the highest posterior probability under the model θ . The criterion for the least confident strategy only considers information about the most

probable label. We can use max margin sampling to preserve information about the remaining label distribution:

$$x_M^* = \arg \min_x P_\theta(\hat{y}_1|x) - P_\theta(\hat{y}_2|x) \quad (141)$$

where \hat{y}_1 and \hat{y}_2 are the first and second most probable class labels under the model, respectively. Intuitively, instances with large margins are easy to classify. Thus, points with small margins are ambiguous and knowing their labels would help the model to more effectively discriminate between them. However, for multi-class problems with very large label sets, the margin strategy still ignores much of the output distribution for the remaining classes. A more general uncertainty sampling strategy is based on entropy:

$$x_H^* = \arg \max_x - \sum_i P_\theta(y_i|x) \log P_\theta(y_i|x) \quad (142)$$

By learning labels that have highest entropy we can reduce label uncertainty. Uncertainty sampling also works for regression problems, in which case the learner queries the point with highest output variance in its prediction.

Query by Committee. Another query selection framework is the Query By Committee (QBC) algorithm that involves maintaining a committee $C = \{\theta^{(1)}, \dots, \theta^{(C)}\}$ of models which are all trained on the current labeled set L but represent competing hypotheses. Each committee member is then allowed to vote on the labelings of query candidates and the most informative query is considered to be an instance about which they most disagree. The objective of QBC is to minimize a set of hypotheses that are consistent with the current labeled training data L . For measuring the level of disagreement two main approaches have been proposed [35]: the vote entropy and KL divergence. The vote entropy is defined as follows:

$$x_{VE}^* = \arg \max_x - \sum_i \frac{V(y_i)}{C} \log \frac{V(y_i)}{C} \quad (143)$$

where $y_i \in \{1, \dots, K\}$ is the class label, $V(y_i)$ is the number of votes that a label received from the committee members and C is the size of the committee. The KL divergence for QBC voting is defined as follows:

$$x_{KL}^* = \arg \max_x \frac{1}{C} \sum_{c=1}^C KL(P_{\theta^{(c)}} || P_C) \quad (144)$$

$$KL(P_{\theta^{(c)}} || P_C) = \sum_i P_{\theta^{(c)}}(y_i|x) \log \frac{P_{\theta^{(c)}}(y_i|x)}{P_C(y_i|x)} \quad (145)$$

$$P_C(y_i|x) = \frac{1}{C} \sum_{c=1}^C P_{\theta^{(c)}}(y_i|x) \quad (146)$$

where $\theta^{(c)}$ represents a member model of the committee and $P_C(y_i|x)$ is the consensus probability that y_i is the correct label. The KL divergence metric considers the most informative query to be the one with the largest average difference between the label distributions of any one committee member and the consensus distribution.

Variance Reduction. We can reduce the generalization error by minimizing output variance. Consider a regression problem for which the learning objective is to minimize the mean squared

error. Let $\bar{\theta} = E[\hat{\theta}]$ be the expected value of the parameter estimate $\hat{\theta}$ and let θ^* be the ground truth, then

$$\text{MSE} = E[(\hat{\theta} - \theta^*)^2] = E[(\hat{\theta} - \bar{\theta}) + (\bar{\theta} - \theta^*)]^2 \quad (147)$$

$$= E[(\hat{\theta} - \bar{\theta})^2] + 2(\bar{\theta} - \theta^*)E[\hat{\theta} - \bar{\theta}] + (\bar{\theta} - \theta^*)^2 \quad (148)$$

$$= E[(\hat{\theta} - \bar{\theta})^2] + (\bar{\theta} - \theta^*)^2 \quad (149)$$

$$= \text{VAR}[\hat{\theta}] + \text{bias}^2(\hat{\theta}) \quad (150)$$

This is called the **bias-variance tradeoff**. Thus, it is possible to achieve lower MSE with a biased estimator as long as it reduces the variance. A natural question is how low can the variance be? The answer is given by the Cramer-Rao lower bound that provides a lower bound on the variance of any unbiased estimator.

Theorem 3.1 (*Cramer-Rao Lower Bound*) Assuming $p(x|\theta)$ satisfies the regularity condition, the variance of any unbiased estimator $\hat{\theta}$ satisfies:

$$\text{VAR}(\hat{\theta}) \geq \frac{1}{-E\left[\frac{\partial^2 \log p(x|\theta)}{\partial \theta^2}\right]} = \frac{1}{I(\theta)} \quad (151)$$

where $I(\theta)$ is the Fisher information matrix.

Thus, the Minimum Variance Unbiased (MVU) estimator achieves the minimum variance equal to the inverse of the Fisher information matrix. To minimize variance of parameter estimates, an active learner should select data that maximizes its Fisher information. For multi-variate models with K parameters, Fisher information takes the form of a $K \times K$ matrix:

$$[I(\theta)]_{ij} = -E\left[\frac{\partial^2 \log p(x|\theta)}{\partial \theta_i \partial \theta_j}\right] \quad (152)$$

As a result, there are several options for minimizing the inverse information matrix: A-optimality minimizes the trace: $\text{Tr}(I^{-1}(\theta))$, D-optimality minimizes the determinant: $|I^{-1}(\theta)|$ and E-optimality minimizes the maximum eigenvalue: $\lambda_{\max}[I^{-1}(\theta)]$.

However, there are some computational disadvantages to the variance reduction methods. Estimating output variance requires inverting a $K \times K$ matrix for each unlabeled instance, resulting in a time complexity of $\mathcal{O}(UK^3)$, where U is the size of the query pool. As a result variance reduction methods are empirically slower than simpler query strategies like uncertainty sampling.

Active learning and *semi-supervised learning* both try to make the most of unlabeled data. For example, a basic semi-supervised technique is self-training in which the learner is first trained with a small amount of labeled data and then used to classify the unlabeled data. The most confident unlabeled instances together with their predicted labels are added to the training set and the process repeats.

Figure 33 compares 3 uncertainty sampling techniques: least confident, max margin and entropy with a random subsampling baseline on 20 newsgroups dataset classified with Logistic

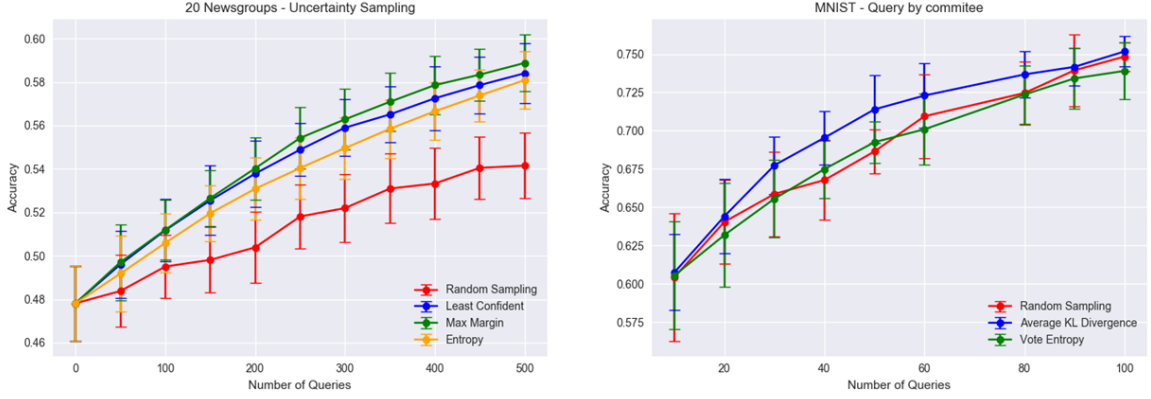


Figure 33: Active Learning: Uncertainty Sampling and Query By Committee.

Regression. All three methods achieve higher accuracy in comparison to baseline that highlights the benefit of active learning. On the right, we can see the performance of Query By Committee strategy applied to MNIST dataset. The committee consists of 5 instances of Logistic Regression. Two methods are compared against the random subsampling baseline: vote entropy and average KL divergence. We can see that average KL divergence achieves highest classification accuracy. All experiments were repeated 10 times.

3.4 Submodular Optimization

Many observation selection problems satisfy an intuitive **diminishing returns** property: adding an observation helps more if we made few observations so far, and helps less if we already made a lot of observations. From a sensing perspective, combining two sets of information has less gain than the sum of information gains of the individual sets. This concept is formalized by the property of submodularity.

Submodularity is a property of *set functions* $f : 2^V \rightarrow \mathbb{R}$ that assign each subset $S \subseteq V$ a value $f(S)$. F is called *submodular* if for all $A, B \subseteq V$ it holds that

$$F(A) + F(B) \geq F(A \cup B) + F(A \cap B) \quad (153)$$

Note that if F is non-negative ($F(A) \geq 0$ for all A) submodularity implies subadditivity: $f(A \cup B) \leq f(A) + f(B)$. An alternative definition of submodular functions is the following. A set function is sub-modular iff for all $A \subseteq B \subseteq V$ and $s \in V \setminus B$, it holds that

$$F(A \cup \{s\}) - F(A) \geq F(B \cup \{s\}) - F(B) \quad (154)$$

We can usually check more easily for monotonicity and submodularity of a function via its incremental definition, where the increment function $f(A|B)$ is defined as $f(A|B) = f(A \cup B) - f(B)$. Thus, for a single element $\{e\} \in V \setminus B$ and $A \subseteq B \subseteq V$, we can define sub-modularity as:

$$f(e|A) \geq f(e|B) \quad (155)$$

An important subclass of submodular functions are those which are *monotone* where enlarging the argument set cannot cause the function to decrease. A function $f : 2^V \rightarrow \mathbb{R}$ is monotone if for every $A \subseteq B \subseteq V$, we have $f(A) \leq f(B)$.

Submodularity is closed under the operations of non-negative addition, restriction and conditioning.

Theorem 3.2 (*Submodularity is closed under conic combination*). Given submodular functions f_1, f_2, \dots, f_n and $\alpha_i \geq 0$, their conic combination (non-negative addition) is submodular.

$$f(A) = \sum_{i=1}^n \alpha_i f_i(A) \quad (156)$$

Proof. Since every function f_i is submodular and $\alpha_i \geq 0$, we have for $A \subseteq B$:

$$f_i(j|A) \geq f_i(j|B) \Rightarrow \sum_{i=1}^n \alpha_i f_i(j|A) \geq \sum_{i=1}^n \alpha_i f_i(j|B) \Rightarrow f(j|A) \geq f(j|B) \quad (157)$$

Theorem 3.3 (*Entropy is submodular*). Let V be the index set of a set of random variables, then the following function is submodular:

$$f(A) = H(X_A) = - \sum_{X_A} p(X_A) \log p(X_A) \quad (158)$$

Proof. For $A \subseteq B \subseteq V$ and $e \in V \setminus B$, we have:

$$F(A \cup e) - F(A) = H(X_A, X_e) - H(X_A) = H(X_e|X_A) \quad (159)$$

$$\geq H(X_e|X_A, X_B) = H(X_e|X_B) = H(X_B, X_e) - H(X_B) \quad (160)$$

$$= F(B \cup e) - F(B) \quad (161)$$

where the inequality is due to the fact that conditioning reduces entropy.

Theorem 3.4 (*Mutual Information is submodular*). MI is a submodular function if observed variables are conditionally independent given the latent state.

Proof. MI can be expressed as a set function $f(A) = I(X; Y_A)$. We make a distinction between latent variables X and observed variables Y . For any pair $I, J \subseteq V$ such that $I \cap J = \emptyset$, we assume Y_I is independent of Y_J given X . Let $A \subseteq B \subseteq V$, and $j \in V \setminus B$, we have:

$$I(Y_j; Y_{B \setminus A} | Y_A) \geq 0 \quad (162)$$

$$H(Y_j | Y_A) - H(Y_j | Y_{(B \setminus A) \cup A}) \geq 0 \quad (163)$$

$$H(Y_j | Y_A) \geq H(Y_j | Y_B) \quad (164)$$

which is the submodularity for entropy. Continuing by using conditional independence assumption and subtracting $H(Y_j | X) = H(Y_j | X, Y_A) = H(Y_j | X, Y_B)$ from both sides:

$$H(Y_j | Y_A) - H(Y_j | X) \geq H(Y_j | Y_B) - H(Y_j | X) \quad (165)$$

$$H(Y_j | Y_A) - H(Y_j | X, Y_A) \geq H(Y_j | Y_B) - H(Y_j | X, Y_B) \quad (166)$$

$$I(Y_j; X | Y_A) \geq I(Y_j; X | Y_B) \quad (167)$$

$$I(X; Y_j | Y_A) \geq I(X; Y_j | Y_B) \quad (168)$$

3.4.1 Greedy Maximization

Submodular functions arise in many applications and therefore it is natural to study submodular optimization. In particular, we focus on greedy maximization of submodular functions. We are interested in solving problems of the form:

$$\max_{S \subseteq V} f(S) \quad (169)$$

subject to some constraints on S . We will focus on two kinds of constraints: *cardinality constraint* $|S| \leq k$ and the *budget constraint* $\text{cost}(S) \leq B$. An example of the former is a sensor placement problem where we want to find the k best sensor locations. An example of the latter is a knapsack problem where we would like to maximize the value of items given a total budget of their weight.

A simple greedy approach to maximizing submodular functions is a greedy selection summarized in Algorithm 7 that starts with an empty set S_0 and at iteration i , adds the element maximizing the increment function:

$$S_i = S_{i-1} \cup \{\arg \max_e f(e|S_{i-1})\} \quad (170)$$

Algorithm 7 Greedy Algorithm for Maximizing a Submodular Function

```

1: Init  $A_0 = \emptyset$ ,  $k$ 
2: for  $j = 1, 2, \dots, k$  do
3:    $a_j = \arg \max_{e \in V} f(e|A_{j-1})$ 
4:   set  $A_j = A_{j-1} \cup \{a_j\}$ 
5:   set  $V = V \setminus \{a_j\}$ 
6: end for

```

A celebrated result by Nemhauser [27] proves that the greedy algorithm provides a good approximation to the optimal solution of the NP-hard optimization problem.

Theorem 3.5 *Let $f : 2^V \rightarrow \mathbb{R}$ be a non-negative, monotone submodular function and let $\{S_i\}_{i \geq 0}$ be the greedily selected sets. Then for all positive integers k and l :*

$$f(S_l) \geq \left(1 - e^{-l/k}\right) \max_{S: |S| \leq k} f(S) \quad (171)$$

Proof. Let $S^* = \arg \max\{f(S) : |S| \leq k\}$ be an optimal set of size k . Order the elements of S^* arbitrarily as $\{v_1^*, \dots, v_k^*\}$. Then, we have the following sequence of inequalities for all $i < l$:

$$f(S^*) \leq f(S^* \cup S_i) \quad (172)$$

$$= f(S_i) + \sum_{j=1}^k f(v_j^* | S_i \cup \{v_1^*, \dots, v_{j-1}^*\}) \quad (173)$$

$$\leq f(S_i) + \sum_{v \in S^*} f(v | S_i) \quad (174)$$

$$\leq f(S_i) + \sum_{v \in S^*} (f(S_{i+1}) - f(S_i)) \quad (175)$$

$$\leq f(S_i) + k(f(S_{i+1}) - f(S_i)) \quad (176)$$

where the first inequality follows from monotonicity of f , the second equality is a telescoping sum, the third inequality follows from the submodularity of f , the fourth holds because S_{i+1} is built greedily from S_i , and the last inequality reflects the fact that $|S^*| \leq k$. Therefore,

$$f(S^*) - f(S_i) \leq k(f(S_{i+1}) - f(S_i)) \quad (177)$$

Define $\delta_i = f(S^*) - f(S_i)$, we can then re-write the above expression as $\delta_i \leq k(\delta_i - \delta_{i+1})$, that can be re-arranged to yield:

$$\delta_{i+1} \leq \left(1 - \frac{1}{k}\right) \delta_i \quad (178)$$

Therefore, $\delta_l \leq (1 - \frac{1}{k})^l \delta_0$. Using the inequality $1 - x \leq e^{-x}$, we have

$$\delta_l \leq \left(1 - \frac{1}{k}\right)^l \delta_0 \leq e^{-l/k} f(S^*) \quad (179)$$

where we used the fact that $\delta_0 = f(S^*) - f(\emptyset) \leq f(S^*)$ since f is non-negative by assumption. Substituting $\delta_l = f(S^*) - f(S_l)$ and re-arranging, we get:

$$f(S_l) \geq (1 - e^{-l/k}) f(S^*) \quad (180)$$

Therefore, if we let the greedy algorithm pick $l = k$ sensors (compared to the optimal set of size k), the greedy solution is no worse than $1 - 1/e \approx 0.632$ of the optimal value!

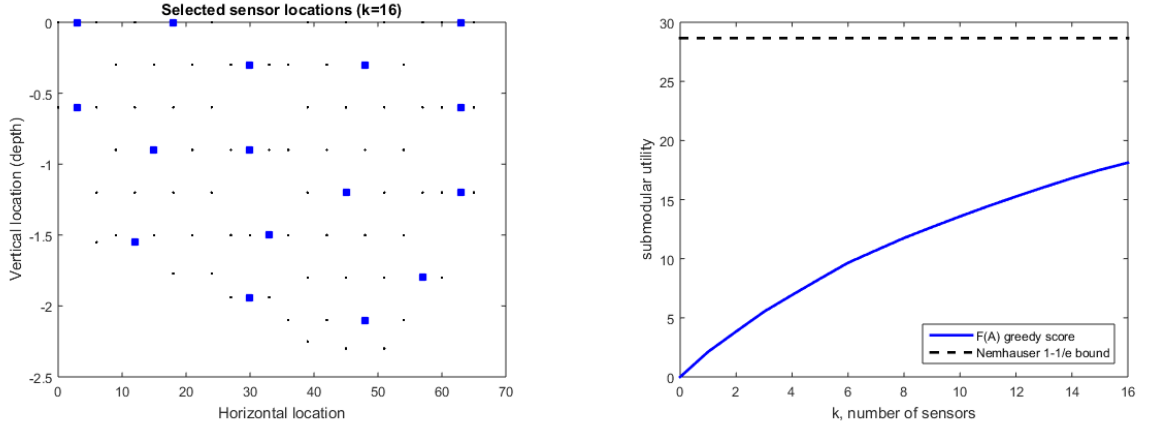


Figure 34: Submodular function optimization subject to cardinality constraint.

Figure 34 shows a sensor selection problem that maximizes mutual information $F_{mi}(A) = H(V \setminus A) - H(V \setminus A|A)$ over the ground set V subject to cardinality constraint: $|S| \leq k$. We can see that the selected sensors evenly cover the geographic area. Also shown are submodular utility scores along with Nemhauser bound obtained by taking the last score and dividing it by $(1 - 1/e)$.

In many applications, the elements $v \in V$ may have non-uniform costs $c(v) \geq 0$ and we may want to maximize f subject to a budget B that the total cost cannot exceed:

$$\max_S f(S) \text{ s.t. } \sum_{v \in S} c(v) \leq B \quad (181)$$

The standard (uniform cost) greedy algorithm can perform arbitrarily badly since it ignores the cost. Therefore, we need to modify the selection step by normalizing the increment function by the element's cost:

$$S_{i+1} = S_i \cup \left\{ \arg \max_{v \in V \setminus S_i : c(v) \leq B - c(S_i)} \frac{f(v|S_i)}{c(v)} \right\} \quad (182)$$

Thus, we choose an element that has highest incremental gain and lowest cost. Note that the sensor selection problem can be thought of as a special case of the budget constrained problem if we assign unit costs to all measurements and set the budget $B = k$. Let S_{uc} be the uniform cost solution and S_{cb} be the cost-benefit greedy solution described above. It can be shown that:

$$\max\{f(S_{uc}), f(S_{cb})\} \geq \frac{1 - 1/e}{2} \text{OPT} \quad (183)$$

The submodular maximization algorithm with a budget constraint also known as Cost-Effective Lazy Forward-selection (CELf) Algorithm is summarized in Algorithm 8.

Algorithm 8 Cost-Effective Lazy Forward-selection (CELF)

```
1: Input: Reward function  $F$ , cost function  $C$ , budget  $B$ 
2:    $A_{UC} = \text{LazyForward}(F, C, B, 'UC')$ 
3:    $A_{CB} = \text{LazyForward}(F, C, B, 'CB')$ 
4:   return  $\arg \max\{F(A_{UC}), F(A_{CB})\}$ 
5:  $\text{LazyForward}(F, C, B, \text{type})$ 
6:
7: end for
```

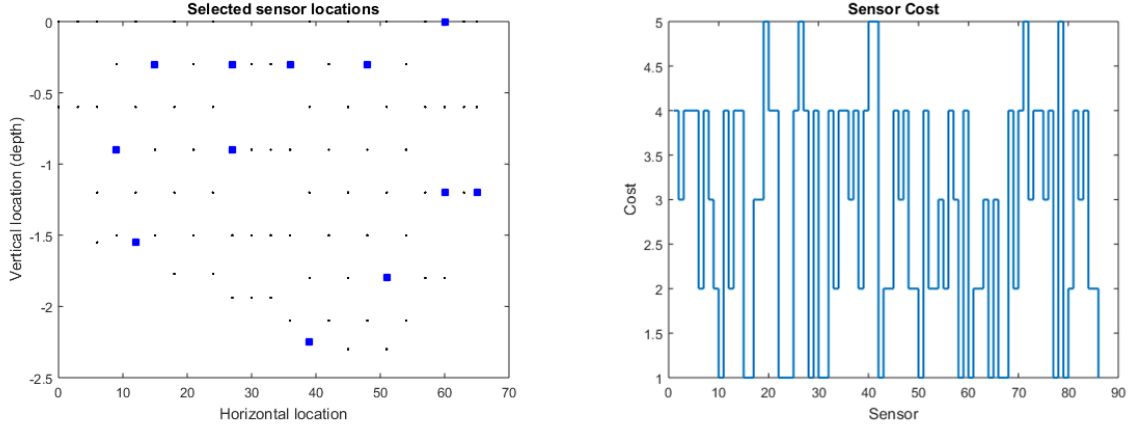


Figure 35: Submodular function optimization subject to budget constraint.

4 Information Planning

Information planning enables faster learning with fewer training examples. It is particularly applicable when training examples are costly to obtain. This work examines the advantages of information planning for text data by focusing on three supervised models: Naive Bayes, supervised LDA and deep neural networks. We show that planning based on entropy and mutual information outperforms random selection baseline and therefore accelerates learning.

4.1 Introduction

Information planning involves making decisions based on information measures [7]. Information planning is closely related to active learning [35], [28] and optimum experiment design [10], [6] in which labeled data is expensive to obtain. The key idea behind information planning is that a model can learn better with fewer training examples if the training examples are carefully selected to maximize the information gain for a particular task. In other words, given an abundance of unlabelled data, we would like to rank the unlabelled examples according to their usefulness in training of the model. The top K most informative training examples are labelled by expert annotators and added to the training set.

There are three main scenarios in which the learner may be able to ask queries: membership queries, stream-based selective sampling and pool-based sampling. In membership query, the learner may request labels for any unlabelled instance in the input space. In stream-based selective sampling, the active learning is done sequentially as each unlabelled instance is drawn one at a time from the data source and the learner must decide whether to query or discard it. Finally, in pool-based sampling, the entire set of unlabelled data is ranked and the most informative subset is labelled and added to the training set. In this work, we focus on pool

Algorithm 9 Generic Information Planning

```
1: Input: data  $x \in \{D_{train}, D_{pool}\}$ , model  $M$ , acquisition function  $a(x, M)$ 
2: Output: trained model  $M^{(T)}$ 
3: for  $t = 1, 2, \dots, T$  do
4:   train  $M^{(t)}$  on  $D_{train}$ 
5:   select  $x^* = \arg \max_{x \in D_{pool}} a(x, M^{(t)})$ 
6:   update  $D_{train} = D_{train} \cup x^*$ ;  $D_{pool} = D_{pool} \setminus x^*$ 
7: end for
```

based sampling as summarized in Algorithm 9.

The decisions whether or not to query instances of unlabelled examples are based on a measure of informativeness or a query strategy. The simplest and most commonly used query framework is *uncertainty sampling* [22]. In this framework, an active learner queries instances about which it is least certain about the label. For example, in a binary classification task, uncertainty sampling queries points near the decision boundary where the class probabilities are close to $1/2$, i.e. highest entropy. Uncertainty sampling also works for regression problems, in which case the learner queries points with highest output variance in their prediction. In general, uncertainty sampling chooses a subset x^* from a set of unlabelled examples $x \in D_{pool}$ that maximizes label entropy:

$$x^* = \arg \max_{x \in D_{pool}} - \sum_c p(y = c|x) \log p(y = c|x) \quad (184)$$

An alternative measure that not only looks at the uncertainty of the label but also at the informativeness is mutual information:

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) \quad (185)$$

For random variables X and Y , to maximize the *information gain* $I(X; Y)$, we want to maximize $H(X)$ (uncertainty about X) and minimize $H(X|Y)$, i.e. we want Y to be informative about X . The choice of random variables X and Y depends on the specific application and probabilistic model. As we will see, in complex models, often one direction of mutual information is easier to compute compared to the other.

Other query selection strategies include ensemble methods such as query-by-committee [36] in which an ensemble of diverse and accurate models are trained on the labeled dataset. Each ensemble member is then allowed to vote on the labels of query candidates and the most informative query is considered to be an instance about which they most disagree. In this work, we focus on uncertainty sampling and mutual information as the primary query strategies.

In application to natural language processing, active learning has been successfully applied to information extraction [33], named entity recognition [1], text categorization [41], part-of-speech tagging [32] and parsing [2]. In this work, we study the benefits of information planning for text data on three different tasks: text classification, supervised topic modeling and sentiment classification.¹ To demonstrate the advantage of information planning, we compare performance of active learning against random selection. As the size of unlabelled dataset increases, we expect active learning to perform better.

4.2 Naive Bayes

We consider the task of text classification using a Naive Bayes graphical model as shown in Figure 36.

¹all code is available on-line

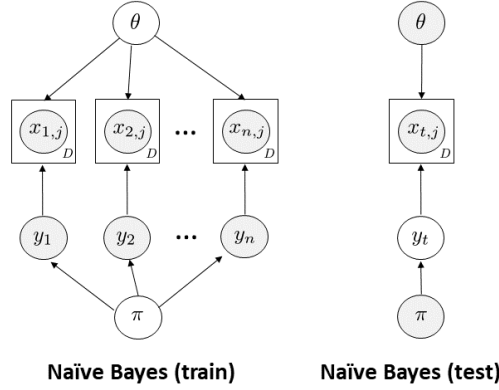


Figure 36: Naïve Bayes graphical model.

Let x_{ij} be Bernoulli random variables indicating the presence ($x_{ij} = 1$) or absence ($x_{ij} = 0$) of a word $j \in \{1, \dots, D\}$ in document $i \in \{1, \dots, n\}$, parameterized by θ_{jc} for a given class label $y_i \in \{1, \dots, C\}$. In addition, let π be a Dirichlet distribution representing the prior over the class labels. Thus, the total number of parameters is $|\theta| + |\pi| = \mathcal{O}(DC) + \mathcal{O}(C) = \mathcal{O}(DC)$. Due to the small number of parameters, the Naïve Bayes model is immune to over-fitting [26].

The choice of Bernoulli Naïve Bayes formulation is important because it leads itself to word-based information planning. By associating each word in the dictionary with a binary random variable, we are able to compute the influence of individual words on class label distribution.

Consider words x_i in a single document i :

$$p(x_i, y_i | \theta) = p(y_i | \pi) \prod_{j=1}^D p(x_{ij} | y_i, \theta) = \prod_{c=1}^C \pi_c^{1[y_i=c]} \prod_{j=1}^D \prod_{c=1}^C p(x_{ij} | \theta_{jc})^{1[y_i=c]} \quad (186)$$

We can write-down the Naïve Bayes inference algorithm by maximizing the log-likelihood objective:

$$\begin{aligned} \log p(D | \theta) &= \log \prod_{i=1}^n p(x_i, y_i | \theta) = \sum_{i=1}^n \log p(x_i, y_i | \theta) \\ &= \sum_{c=1}^C N_c \log \pi_c + \sum_{j=1}^D \sum_{c=1}^C \sum_{i: y_i=c} \log p(x_{ij} | \theta_{jc}) \end{aligned} \quad (187)$$

During test time, we would like to predict the class label y given the training data D and the learned model parameters. Applying the Bayes rule:

$$p(y = c | x_{i,1}, \dots, x_{i,D}, D) \propto p(y = c | D) \prod_{j=1}^D p(x_{ij} | y = c, D) \quad (188)$$

Substituting $p(x_{ij} | y = c, D)$ and taking the log, we get:

$$\log p(y = c | x, D) \propto \log p(y = c | D) + \sum_{j=1}^D (1[x_{ij} = 1] \log \theta_{jc} + 1[x_{ij} = 0] \log(1 - \theta_{jc})) \quad (189)$$

In entropy based planning we want to rank test documents according to the entropy of their class label. In the case of class distribution, the entropy can be computed in closed form: $H(y) = -\sum_{c=1}^C p(y = c) \log p(y = c)$.

In the case of mutual information, we need to choose the random variables we want to use for planning. Since we are interested in classifying documents, it makes sense to choose y_i (the class label for document i) as one of the variables. Our choice of the second variable depends on the quantity of interest. We can choose one of the global variables θ_{jc} (probability that word j is present in class c) or π_c (probability of class c). In this example, we consider estimating $I(y_i; \theta)$, i.e. we are interested in measuring the information gain about the word distribution θ given the class label y_i for a test document. Since both variables are discrete, the mutual information can be estimated in closed form:

$$I(y_i; \theta) = \text{KL}(p(y_i, \theta) || p(y_i)p(\theta)) = \sum_{y \in C} \sum_{\theta} p(\theta|y_i)p(y_i) \log \frac{p(\theta|y_i)}{p(\theta)} \quad (190)$$

We can compute $p(\theta) = \sum_{c \in C} p(x_j = 1 | y = c)p(y = c)$ for $j \in \{1, \dots, D\}$. In addition, we compute $I(x_j; \pi)$ to measure how informative each word x_j is to the global label distribution π .

We train our Naive Bayes classifier on a subset of the 20newsgroups dataset. In particular, we'll restrict ourselves to 4 classes: space, graphics, autos and hockey. We'll use a count vectorizer to produce a vector of word counts for each document while filtering stop and low frequency words. Figure 37 shows the posterior word distribution for the four classes.

Let's visualize the test document ranking based on entropy and mutual information in Figure 38. In particular, we are interested in outliers, i.e. documents that have highest entropy and mutual information. To visualize the difference between entropy and MI based planning, we can sort the documents according to entropy and use the same index to sort MI as shown in Figure 39. In Figure 39 documents that have high entropy but low MI are uncertain but not informative. Therefore, we expect MI to provide a better measure for information planning.

To evaluate the utility of information planning in comparison to random selection, we take a total of 1K documents and divide them into 100 labelled training documents, 700 unlabelled test documents, and holdout 200 documents for evaluating classification performance. We run the active learning pipeline, in which top K most informative documents are selected from the test set, labelled by annotators and added to the training set, at which point the model is re-trained from scratch and evaluated on the 200 heldout documents. Figure 40 shows the experimental results of this active learning pipeline. We can see a clear advantage of using information planning compared to random selection.

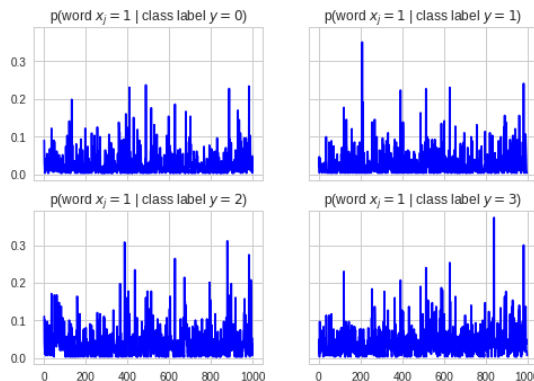


Figure 37: Naive Bayes class-conditional posterior.

4.3 Supervised LDA

Latent Dirichlet Allocation (LDA) [4] is an unsupervised topic model that represents each document as a mixture of topics, where a topic is a distribution over words. The objective is to learn the shared topic distributions and their proportions for each document. However, it is often desirable to associate a quantity of interest (e.g. a score) with each document; this could be a five star rating for a movie review, a quality score of a report, or a number of times an on-line article is linked. Supervised LDA (sLDA) [3] is designed to jointly learn the topics and the score variable associated with each document. The graphical model of supervised LDA is shown in Figure 41.

sLDA associates each word x_{id} with a topic label $z_{id} \in \{1, \dots, K\}$. Each document is associated with topic proportions θ_d that can also be used as a measure of document similarity. The topics β_k represented by the Dirichlet distribution are shared across all documents. The hyper-parameters α and η capture our prior knowledge of topic proportions and the topics, e.g. from past training of the model. Finally, y_d is the response variable that represents a score associated with each document. In this example, y_d is modelled as a regression between global coefficients w and a normalized histogram of topics present in a document: \bar{z} . Assuming a Gaussian distribution, y_d can be expressed as follows:

$$y_d \sim N(w_1 \bar{z}_1 + \dots + w_K \bar{z}_K, \sigma^2) \quad (191)$$

$$\bar{z} = \frac{1}{N} \sum_{n=1}^N z_n \quad (192)$$

where z_n is a one-hot encoded topic assignment vector of size $1 \times K$. We use PyMC [30] Metropolis-Hastings sampler to infer the latent variables. Figure 42 shows the first 4 (out of 8) topics of sLDA model trained on the polarity sentiment dataset [29]. Notice, that the addition of the score variable y_d influences topic inference due to the coupling between w and β_k introduced by the $z_{i,d} \rightarrow y_d \leftarrow w$ and $z_{i,d} \rightarrow x_{i,d} \leftarrow \beta_k$ v-structures [21].

We can use a sample based estimator to compute the entropy of the score variable y_d : $H[y_d|x_d, D] = E[-\log p(y_d|x_d, D)]$.

Assuming we have topic samples:

$$\{(z_{1d}^{(i)}, \dots, z_{Nd}^{(i)})\}_{i=1}^M \overset{MCMC}{\sim} p(z_{1d}, \dots, z_{Nd}|x_d, D) \quad (193)$$

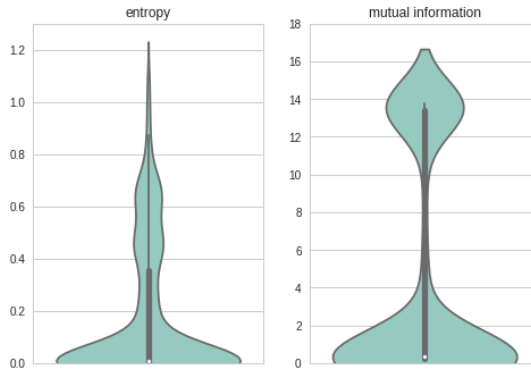


Figure 38: Violin plot showing the density of entropy and MI on test data.

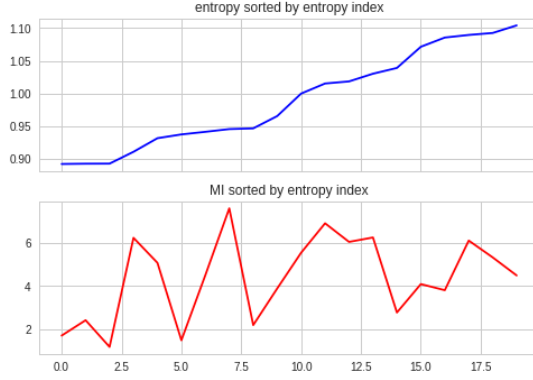


Figure 39: Test document entropy and MI sorted by entropy index.

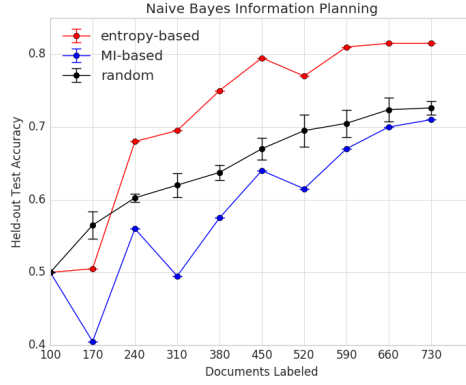


Figure 40: Classification accuracy as a function of labeled examples.

$$p(y_d|x_d, D) = \sum_{z_{1d}} \dots \sum_{z_{Nd}} p(z_{1d}, \dots, z_{Nd}|x_d, D) p(y_d|z_d) \approx \frac{1}{M} \sum_{i=1}^M p(y_d|z_d = z_d^{(i)}) = \hat{p}(y_d|x_d, D) \quad (194)$$

Once we have a sample estimate $\hat{p}(y_d|x_d, D)$, to compute entropy, we need to estimate the expectation:

$$E[-\log \hat{p}(y_d|x_d, D)] = - \int p(y_d|x_d, D) \log \hat{p}(y_d|x_d, D) dy_d \approx \frac{1}{M} \sum_{i=1}^M \log \hat{p}(y_d^{(i)}|x_d, D) \quad (195)$$

where $\{(y_d^{(i)})\}_{i=1}^M \stackrel{MCMC}{\sim} \hat{p}(y_d^{(i)}|x_d, D)$.

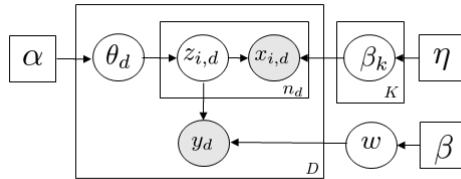


Figure 41: Supervised LDA [3] graphical model.

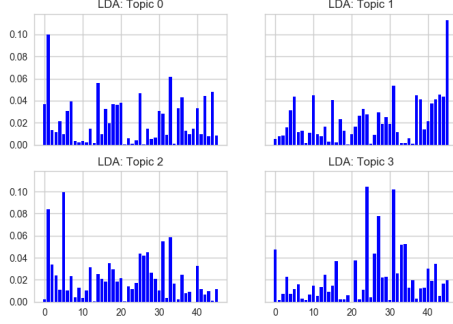


Figure 42: Supervised LDA topics inferred using PyMC [30].

To evaluate the utility of information planning in comparison to random selection, we start with a corpus of 100 documents, from which we allocate 15 documents for training, 35 for testing and holdout 50 documents for evaluating sLDA performance. Figure 43 shows the test MSE of the Gaussian score variable evaluated against the ground truth review score on the heldout set of 50 documents.

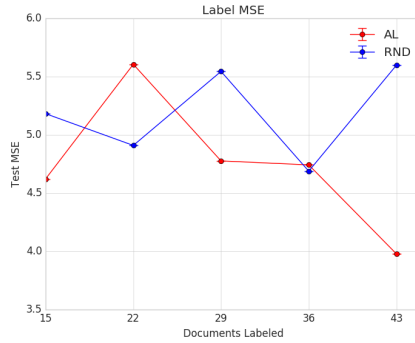


Figure 43: Supervised LDA test MSE.

We can see that active learning selection achieves smaller MSE in comparison to random selection baseline.

4.4 Deep Neural Network

We consider information planning in the case of deep neural networks. In particular, we examine two architectures: LSTM and CNN in application to sentence sentiment classification [15].

Recurrent Neural Networks (RNNs) provide a natural way of encoding a sentence as a sequence of words. Just as you are processing this sentence word by word while keeping a memory of what came before, RNNs maintain an internal model of the past and update it as soon as the new information arrives. Thus, RNNs contain an internal loop unrolled over the length of the input sequence. As a result, RNNs can capture more information in comparison to n -gram models. For example, in a n -gram language model (with $n = 3$), we can predict the next word given the previous two:

$$l_{sent} = \frac{1}{n} \log P(w_1, \dots, w_n) = \frac{1}{n} \sum_{i=1}^n \log p(w_i | w_{i-1}, w_{i-2}) \quad (196)$$

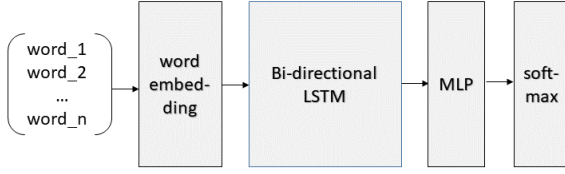


Figure 44: LSTM model architecture

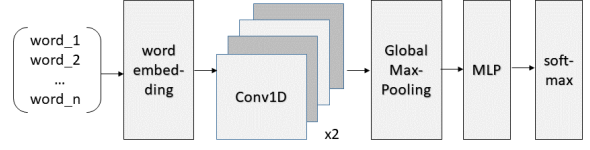


Figure 45: CNN model architecture

The dependence in neural language model (LM) extends to the beginning of the sentence, i.e. $p(w_i|w_{i-1}, \dots, w_1)$, while in a n -gram LM the conditioning is only on the previous n words and as n grows larger, sparsity of training data becomes a problem. In both cases, the word conditional distribution greatly reduces the output space of possible words and enables learning of word representations based on their context [25].

A sentence can be encoded by an RNN (such as LSTM) as either the output of its last layer or via average pooling of its intermediate outputs. In both cases, the word position information is preserved in contrast to averaging of word vectors. Figure 44 shows the LSTM architecture used for sentiment classification. We use a bi-directional LSTM that concatenates sequence representations of forward and reverse LSTMs. This enriches order-sensitive representation of the sentence and leads to a slight increase in accuracy. In addition to output layer dropout, we use recurrent dropout that applies the same dropout mask at every time-step (in contrast to a dropout mask that varies randomly between time-steps that could harm the learning process) [12].

Convolutional Neural Networks (CNNs) offer computational advantages due to their parallel nature and comparable performance with time-series data, especially in cases when recent past is not more informative than the distant past (due to translation invariance of CNNs). The same properties that make convnets excel at computer vision also make them highly relevant for sequence processing. Time can be treated as a spatial dimension like the height or width of a 2D image. For example, 1D convolutions can be used to extract patterns from local 1D patch sub-sequences of the input sequence. Because the same transformation is performed on every patch, a pattern learned at a certain position in the sentence can later be recognized at a different position, making 1D convnets translation invariant. While a word-level conv-net operating on word embeddings with a kernel size of 5 should be able to recognize words or word fragments of length 5 or less, a character-level 1D conv-net is able to learn about word morphology. In addition, we can learn hierarchical patterns by stacking several convolutional layers. Figure 45 shows the CNN architecture used for sentiment classification. It consists of a word embedding layer followed by a stack of two convolutional layers separated by a dropout layer. Finally, the activations across each of the filters are pooled via global max pooling and fed into a Multi-Layer Perceptron (MLP) classifier.

For both LSTM and CNN models, we choose pre-trained word embeddings because in an active learning setting, there's usually not enough training data to learn specialized word embeddings for a given task. In order to preserve the latent structure (that could otherwise be destroyed by large gradients from randomly initialized downstream layers), we configure the embeddings to be non-trainable. Both architectures use Glove pre-trained word embeddings [31] based on word co-occurrence for dense representation of words.

As shown in [13] dropout [38] and other stochastic regularization techniques can be used for approximate inference in deep neural networks. This is accomplished by using dropout at

test-time and computing Monte Carlo integration (referred to as MC dropout):

$$\begin{aligned} p(y = c|x, D) &= \int p(y = c|x, w)p(w|D)dw \\ &\approx \int p(y = c|x, w)q_{\theta}^*(w)dw \approx \frac{1}{T} \sum_{t=1}^T p(y = c|x, \hat{w}_t) \end{aligned} \quad (197)$$

where $\hat{w}_t \sim q_{\theta}^*(w)$ and $q_{\theta}^*(w)$ is a variational distribution that minimizes KL divergence to the true model posterior $p(w|D)$:

$$\arg \min_{\theta} KL(q_{\theta}(w)||p(w|D)) \quad (198)$$

where θ are the variational parameters and w are the neural network weights. Averaging forward passes with MC dropout during test time is equivalent to Monte Carlo integration over a Gaussian process posterior approximation [14].

We consider two acquisition functions: one based on entropy and one based on mutual information. An *acquisition function* $a(x, M)$ determines which measurements to take next, and it's a function of the pool of unlabelled examples $x \in D_{pool}$ and the model M :

$$x^* = \arg \max_{x \in D_{pool}} a(x, M) \quad (199)$$

In the case of entropy, we want to choose data points that maximize:

$$H[y|x, D] = - \sum_c p(y = c|x, D) \log p(y = c|x, D) \quad (200)$$

In the case of mutual information, we want to choose data points that maximize the information gain between the predictions y and the model parameters w :

$$I(y; w|x, D) = H[y|x, D] - E_{p(w|D)}[H[y|w, x, D]] \quad (201)$$

where w in our case are the weights of the neural network and the expectation is computed using MC dropout. Thus, we are seeking data points x for which the model is both marginally most uncertain about y and most confident about the predicted label. Notice, that it's easier to compute MI as written above (as opposed to its alternative version) because the entropies are computed in a discrete, low-dimensional output space. This interpretation allows reasoning about uncertainty in deep learning and the application of Bayesian methods to analyze existing deep learning frameworks.

We evaluate the effectiveness of information planning on Kaggle movie review sentiment dataset [20]. The dataset consists of phrases from the Rotten Tomatoes dataset classified into 5 sentiment labels: negative, somewhat negative, neutral, somewhat positive and positive. Obstacles such as sentence negation, sarcasm, and language ambiguity make this a challenging task. We take the first 10K sentences, set aside 10% for training and holdout 10% for evaluating the performance. At each iteration, we select the top $K = 800$ most informative sentences to be annotated and added to the training set. Figures 50, 51, 52, 53 show active learning experimental results.

4.5 Discussion

In this section we discuss information planning results for each model.

Naive Bayes is the simplest model for studying information planning since all latent variables are discrete and information measures can be computed in closed form. By choosing the Bernoulli

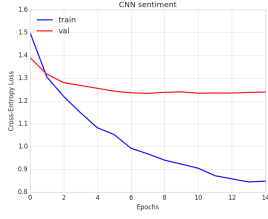


Figure 46: CNN: training loss

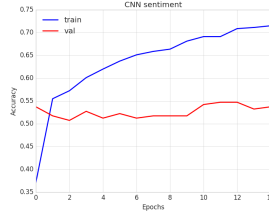


Figure 47: CNN: training acc

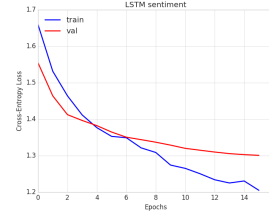


Figure 48: LSTM: training loss

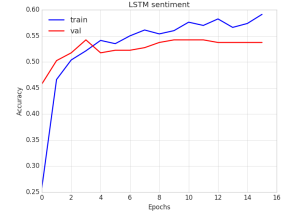


Figure 49: LSTM: training acc

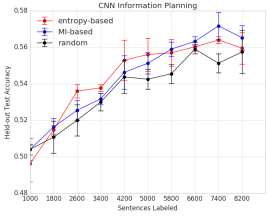


Figure 50: CNN: test accuracy

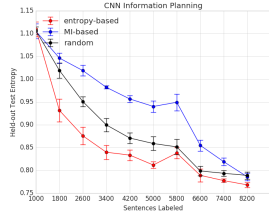


Figure 51: CNN: test entropy

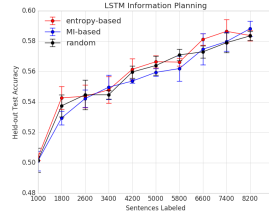


Figure 52: LSTM: test acc.

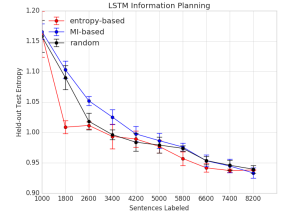


Figure 53: LSTM: test entropy

formulation of Naive Bayes we were able to do planning at the word level. In particular, the top-10 highest MI words as shown in Table 1 were found to be informative of the ground truth labels.

From Figure 40, we can see that we achieve about 10% accuracy increase when we use uncertainty sampling in comparison to random selection. All experimentes were repeated 10 times.

In supervised LDA model, the score variable is learned jointly with the topics. Figure 43 shows that smaller MSE is achieved with uncertainty sampling in comparison to random selection. We also note that sampling based inference for sLDA is more suitable to stream queres in which documents are queried one at a time due considerable time it takes to run sampling for the entire pool of unlabelled examples.

In deep neural network setting we assume that computational expense of re-training the model from scratch is justified by the cost of obtaining the labels. From the training and validation loss in Figures 46, 48, we see no signs of over-fitting due regularization with dropout and weight decay. The validation loss plateaus early in the training indicating that model capacity can be increased in an attempt to achieve higher classification accuracy. Figures 47, 49 show that CNN and LSTM model achieve comparable accuracy on the sentiment classification task with validation accuracy close to 55%.

Figures 50, 52 show the benefit of information planning over the random baseline. We can see that we achieve close to 1% improvement in accuracy in the case of CNN and a relatively smaller improvement in the case of LSTM, even though LSTM achieves 2% higher held-out test accuracy in the end. It's also interesting to see, the decrease in test label entropy as shown in Figures 51, 53, where as expected the uncertainty sampling leads to highest reduction in entropy.

class 0	class 1	class 2	class 3	MI
get	think	first	nasa	teams
please	well	time	could	program
use	know	season	much	league
like	also	year	know	moon
one	cars	play	get	playoffs
anyone	get	hockey	also	orbit
thanks	like	one	like	players
graphics	one	would	one	earth
would	would	game	would	games
know	car	team	space	nasa

Table 1: Naive Bayes: top 10 words. Ground truth labels: graphics, autos, hockey, space.

5 Misc

5.1 Sequential Monte Carlo

Particle Filters (PF) is a Sequential Monte Carlo (SMC) method for estimating the internal states of a Switching Linear Dynamic System (SLDS). A set of particles is used to represent the posterior distribution of the states of the Markov process. At each iteration, the particles get re-samples and the ones that explain the observations best survive to the next iteration.

The Gaussian SLDS can be described as follows [8]:

$$z_t \sim P(z_t|z_{t-1}) \quad (202)$$

$$x_t = A(z_t)x_{t-1} + B(z_t)w_t + F(z_t)u_t \quad (203)$$

$$y_t = C(z_t)x_t + D(z_t)v_t + G(z_t)u_t \quad (204)$$

where $y_t \in \mathbb{R}$ denotes the observations, $x_t \in \mathbb{R}$ denotes latent Gaussian states, $z_t \in \mathbb{Z}$ denotes latent discrete states and u_t is a known control signal. The noise processes are iid Gaussian with $w_t \sim N(0, I)$ and $v_t \sim N(0, I)$. This model implies the continuous densities:

$$p(x_t|z_t, x_{t-1}) \sim N(A(z_t)x_{t-1} + F(z_t)u_t, B(z_t)B(z_t)^T) \quad (205)$$

$$p(y_t|x_t, z_t) \sim N(C(z_t)x_t + G(z_t)u_t, D(z_t)D(z_t)^T) \quad (206)$$

along with initial states $x_0 \sim N(\mu_0, \Sigma_0)$ and $z_0 \sim P(z_0)$. The aim of the analysis is to compute the marginal posterior distribution of the discrete states $p(z_{0:t}, y_{1:t})$. The posterior density can be factorized as follows:

$$p(x_{1:t}, z_{1:t}|y_{1:t}) = p(x_{1:t}|y_{1:t}, z_{1:t})p(z_{1:t}|y_{1:t}) \quad (207)$$

where the density $p(x_{1:t}|y_{1:t}, z_{1:t})$ is Gaussian and can be computed analytically if we know the marginal posterior density $p(z_{1:t}|y_{1:t})$. We can re-write the marginal posterior recursively:

$$p(z_{1:t}|y_{1:t}) = \frac{p(y_t|z_{1:t}, y_{1:t-1})p(z_{1:t}|y_{1:t-1})}{p(y_t|y_{1:t-1})} \quad (208)$$

$$= \frac{p(y_t|z_t)p(z_t|z_{1:t-1}, y_{1:t-1})p(z_{1:t-1}|y_{1:t-1})}{p(y_t|y_{1:t-1})} \quad (209)$$

$$\propto p(y_t|z_t)p(z_t|z_{t-1})p(z_{1:t-1}|y_{1:t-1}) \quad (210)$$

which depends only on the current conditional distributions and previously computed $p(z_{1:t-1}|y_{1:t-1})$. The basic idea is to approximate the belief state of the entire state trajectory using a weighted set of particles:

$$p(z_{1:t}|y_{1:t}) \approx \sum_{s=1}^S \hat{w}_t^s \delta_{z_{1:t}^s}(z_{1:t}) \quad (211)$$

We update this belief using importance sampling. If the proposal has the form $q(z_{1:t}^s|y_{1:t})$ then the importance weights are given by:

$$w_t^s \propto \frac{p(z_{1:t}^s|y_{1:t})}{q(z_{1:t}^s|y_{1:t})} \propto \frac{p(y_t|x_t, z_t)p(x_t, z_t|x_{t-1}, z_{t-1})}{q(x_t, z_t|x_{t-1}, z_{t-1})} \quad (212)$$

We can choose the transition prior to be the proposal distribution:

$$q(x_t, z_t|x_{t-1}, z_{t-1}) = p(x_t, z_t|x_{t-1}, z_{t-1}) = p(x_t|x_{t-1}, z_{t-1})p(z_t|z_{t-1}) \quad (213)$$

then the importance weights are given by the likelihood function: $w_t \propto p(y_t|x_t, z_t)$. The particle filter algorithm is summarized in Algorithm 10.

Algorithm 10 Particle Filter Algorithm [8]

- 1: *Sequential Importance Sampling Step*
 - 2: for $i = 1$ to N do
 - 3: $z_t^i \sim p(z_t|z_{t-1}^i)$
 - 4: $x_t^i \sim p(x_t|x_{t-1}^i, z_t^i)$ as in eq. (205)
 - 5: end for
 - 6: for $i = 1$ to N do
 - 7: $w_t^i \propto p(y_t|x_t^i, z_t^i)$
 - 8: end for
 - 9: $\hat{w}_t^i = \frac{w_t^i}{\sum_{s'} w_t^{s'}}$
 - 10: *Selection Step*
 - 11: Multiply particles with respect to importance weights w_t^i
 - 12: to obtain N particles $\{x_{1:t}^i, z_{1:t}^i\}_{i=1}^N$
 - 13: end for
-

The selection step modifies the weighted approximate density p_N to an unweighted density \hat{p}_N by eliminating particles with low importance weights and by multiplying particles with high importance weights. Thus, $p_N(x) = \sum_{i=1}^N w_i \delta(x - x_i)$ is replaced by

$$\hat{p}_N(x) = \sum_{k=1}^N \frac{1}{N} \delta(x - x_k^*) = \sum_{i=1}^N \frac{n_i}{N} \delta(x - x_i) \quad (214)$$

There are many resampling schemes such as multinomial, stratified, systematic and residual. All these algorithms are unbiased and can be implemented in $O(N)$ time.

The Rao-Blackwellized Particle Filter (RBPF) is similar to the PF but we only sample the discrete states. Then for each sample of the discrete states, we update the mean and covariance of the continuous states using Kalman filter updates. In particular, we sample z_t^i and then

propagate the mean μ_t^i and covariance Σ_t^i of x_t with a Kalman filter as follows:

$$\mu_{t|t-1}^i = A(z_t^i)\mu_{t-1|t-1}^i + F(z_t^i)u_t \quad (215)$$

$$\Sigma_{t|t-1}^i = A(z_t^i)\Sigma_{t-1|t-1}^i A(z_t^i)^T + D(z_t^i)D(z_t^i)^T \quad (216)$$

$$S_t^i = C(z_t^i)\Sigma_{t|t-1}^i C(z_t^i)^T + D(z_t^i)D(z_t^i)^T \quad (217)$$

$$y_{t|t-1}^i = C(z_t^i)\mu_{t|t-1}^i + G(z_t^i)u_t \quad (218)$$

$$\mu_{t|t}^i = \mu_{t|t-1}^i + \Sigma_{t|t-1}^i C(z_t^i)^T S_t^{i-1} (y_t - y_{t|t-1}^i) \quad (219)$$

$$\Sigma_{t|t}^i = \Sigma_{t|t-1}^i - \Sigma_{t|t-1}^i C(z_t^i)^T S_t^{i-1} C(z_t^i) \Sigma_{t|t-1}^i \quad (220)$$

The RBPF takes slightly longer to compute but results in more accurate predictions. Figure 54 shows the generated SLDS states (left) and the inferred states (middle) by Particle Filter (PF) and the Rao-Blackwellized version (RBPF). We can see that the inferred states closely

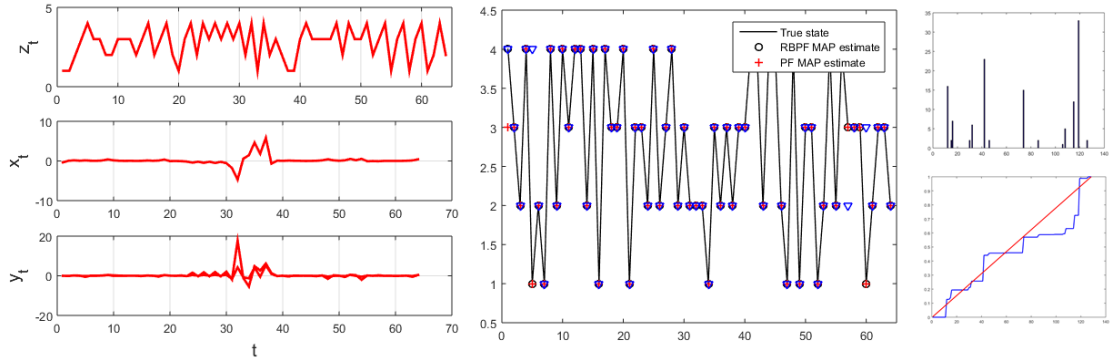


Figure 54: Particle Filter inference using PF and RBPF applied to Switching Linear Dynamic System.

correspond to the ground truth. Also shown is a particle resampling step (right) where only a fraction of the particles survive to the next iteration.

5.2 Information Theory

Information theory addresses the problems of data representation and reliable transmission of information. In order to measure information content, we define several information measures below.

5.2.1 Entropy

The entropy of a random variable is a measure of its uncertainty. Let X be a discrete random variable with probability mass function $p(x)$, then its entropy is defined as:

$$H(X) = -\sum_x p(x) \log p(x) = -E[\log p(x)] \quad (221)$$

A maximum entropy discrete distribution is the uniform distribution. For a random variable with support K and $p(x) = 1/K$, the entropy is $H(X) = -\sum \frac{1}{K} \log \frac{1}{K} = \log_2 K$. Thus, entropy is measured in bits (using log base 2) or nats (using log base e). In contrast, the entropy of a deterministic random variable with $p(x) = \delta[x]$ is $H(X) = 1 \log 1 = 0$. In the special case of a binary random variable $X \in \{0, 1\}$ with $p(x=1) = p$, the entropy is

$$H(X) = -p \log p - (1-p) \log(1-p) \quad (222)$$

It is a concave function with maximum uncertainty of 1 bit occuring at $p = \frac{1}{2}$, when $p = 0$ or $p = 1$, the entropy is 0. Since entropy is a concave function on the space of distributions, we have:

$$H(\lambda p_1 + (1 - \lambda)p_2) \geq \lambda H(p_1) + (1 - \lambda)H(p_2) \quad (223)$$

Since $0 \leq p(x) \leq 1$, we have $\log \frac{1}{p(x)} \geq 0$, and for a discrete X , the entropy is non-negative: $H(X) \geq 0$. We can find an upper bound for entropy using the Jensen's inequality:

$$H(X) = \sum_{x \in \mathcal{X}} p(x) \log \frac{1}{p(x)} \quad (224)$$

$$\leq \log \sum_{x \in \mathcal{X}} p(x) \times \frac{1}{p(x)} \quad (225)$$

$$= \log \sum_{x \in \mathcal{X}} 1 \quad (226)$$

$$= \log |\mathcal{X}| \quad (227)$$

We also note that entropy is independent of permutation or cyclical shifts of the support of our distribution, i.e. it only depends on the point masses $p(x)$.

The joint entropy of a pair of discrete random variables X and Y is defined as:

$$H(X, Y) = - \sum_x \sum_y p(x, y) \log p(x, y) = -E[\log p(x, y)] \quad (228)$$

when two variables are independent the joint entropy is additive: $H(X, Y) = H(X) + H(Y)$ iff $p(x, y) = p(x)p(y)$. The conditional entropy is defined as:

$$H(X|Y) = \sum_y p(y) H(X|Y = y) = - \sum_y p(y) \sum_x p(x|y) \log p(x|y) \quad (229)$$

$$= - \sum_x \sum_y p(x, y) \log p(x|y) = -E[\log p(x|y)] \quad (230)$$

Note that conditioning on Y the uncertainty over X reduces on average: $H(X|Y) \leq H(X)$. The entropy of a pair of variables follows the chain rule:

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y) \quad (231)$$

which follows from the chain rule for probability: $p(x, y) = p(x)p(y|x) = p(y)p(x|y)$. The chain rule can be generalized for multiple random variables X_1, \dots, X_N :

$$H(X_1, \dots, X_N) = \sum_{i=2}^N H(X_i|X_1, \dots, X_{i-1}) + H(X_1) \leq \sum_{i=1}^N H(X_i) \quad (232)$$

where the inequality follows from the fact that conditioning reduces entropy. For continuous random variables, the multivariate Gaussian is the distribution with maximum differential en-

tropy:

$$h(X_1, \dots, X_n) = \int p(x) \log \frac{1}{p(x)} dx \quad (233)$$

$$= \int p(x) \left[\frac{1}{2} \log(2\pi)^n |\Sigma| + \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right] dx \quad (234)$$

$$= \frac{1}{2} \log(2\pi)^n |\Sigma| + \frac{1}{2} E[(x - \mu)^T \Sigma^{-1} (x - \mu)] \quad (235)$$

$$= \frac{1}{2} \log(2\pi)^n |\Sigma| + \frac{1}{2} \text{Tr}\{\Sigma^{-1} \Sigma\} \quad (236)$$

$$= \frac{1}{2} \log(2\pi)^n |\Sigma| + \frac{1}{2} n \log e = \frac{1}{2} \log [(2\pi e)^n |\Sigma|] \quad (237)$$

In information theory, the analog of the law of large numbers is the Asymptotic Equipartition Property (AEP). The AEP states that:

Theorem 5.1 (AEP) *If X_1, X_2, \dots, X_n are iid $\sim p(x)$ then*

$$-\frac{1}{n} \log p(X_1, X_2, \dots, X_n) \rightarrow H(X) \quad \text{in probability} \quad (238)$$

Proof.

$$-\frac{1}{n} \log p(X_1, X_2, \dots, X_n) = -\frac{1}{n} \sum_i \log p(X_i) \rightarrow -E[\log p(X)] = H(X) \quad (239)$$

Thus, the probability $p(X_1, \dots, X_n)$ assigned to an observed sequence will be close to $2^{-nH(X)}$. This enables us to classify the set of all sequences into a typical set, where the sample entropy is close to the true entropy and the nontypical set that contains all other sequences.

5.2.2 KL divergence

One way to measure the similarity between two probability distributions $p(x)$ and $q(x)$ is the Kullback-Leibler divergence or relative entropy. It is defined as follows:

$$KL(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)} \quad (240)$$

we can re-write it as:

$$KL(p||q) = \sum_x p(x) \log p(x) - \sum_x p(x) \log q(x) = -H(p) + H(p, q) \quad (241)$$

where $H(p, q)$ is the cross-entropy. The regular entropy can be written as $H(p, p)$ and therefore KL divergence can be seen as a penalty of extra bits needed to encode the data due to the fact that we used a distribution $q(x)$ to represent the data instead of the true distribution $p(x)$.

Theorem 5.2 (Information Inequality) $KL(p||q) \geq 0$ with equality iff $p = q$.

Proof.

$$-KL(p||q) = -\sum_x p(x) \log \frac{p(x)}{q(x)} = \sum_x p(x) \log \frac{q(x)}{p(x)} \quad (242)$$

$$\leq \log \sum_x p(x) \frac{q(x)}{p(x)} = \log \sum_x q(x) \quad (243)$$

$$\leq \log 1 = 0 \quad (244)$$

where we used the Jensen's inequality, which states that for any convex function $f(x)$, we have

$$f\left(\sum_{i=1}^n \lambda_i x_i\right) \leq \sum_{i=1}^n \lambda_i f(x_i) \quad (245)$$

where $\lambda_i \geq 0$ and $\sum_{i=1}^n \lambda_i = 1$. However, KL divergence is not a true distance between distributions because it is not symmetric ($KL(p||q) \neq KL(q||p)$) and it does not satisfy the triangle inequality. We can use the information inequality to derive an upper bound on entropy. Let $u(x) = 1/|\mathcal{X}|$ be the uniform distribution, then:

$$0 \leq KL(p||u) = \sum_x p(x) \log \frac{p(x)}{u(x)} \quad (246)$$

$$= \sum_x p(x) \log p(x) - \sum_x p(x) \log u(x) \quad (247)$$

$$= -H(X) + \log |\mathcal{X}| \quad (248)$$

Thus, $H(X) \leq \log |\mathcal{X}|$ with equality iff $p(x) = u(x)$.

Theorem 5.3 $KL(p||q)$ is convex in the pair (p, q) , i.e.

$$KL(\lambda p_1 + (1 - \lambda)p_2 || \lambda q_1 + (1 - \lambda)q_2) \leq \lambda KL(p_1 || q_1) + (1 - \lambda)KL(p_2 || q_2) \quad (249)$$

Proof. Expanding the inequality above, we want to show that

$$\sum_x (\lambda p_1(x) + (1 - \lambda)p_2(x)) \log \frac{\lambda p_1(x) + (1 - \lambda)p_2(x)}{\lambda q_1(x) + (1 - \lambda)q_2(x)} \leq \quad (250)$$

$$\lambda \sum_x p_1(x) \log \frac{p_1(x)}{q_1(x)} + (1 - \lambda) \sum_x p_2(x) \log \frac{p_2(x)}{q_2(x)} \quad (251)$$

Then for a fixed x , it suffices to show that

$$\sum_{i=1,2} \lambda_i p_i \log \frac{\lambda_i p_i}{\lambda_i q_i} \geq \left(\sum_{i=1,2} \lambda_i p_i \right) \left(\log \frac{\sum_i \lambda_i p_i}{\sum_i \lambda_i q_i} \right) \quad (252)$$

where $\lambda_1 = \lambda$ and $\lambda_2 = 1 - \lambda$, which is exactly the log-sum inequality.

In order to find a distribution $q(x) \in Q$ that is closest to $p(x) \in P$, we can minimize $KL(q||p)$ with respect to $q(x)$ known as I-projection or information projection:

$$(\text{I} - \text{projection}) : \quad q(x) = \arg \min_{q \in Q} KL(q||p) \quad (253)$$

Since the KL divergence is not symmetric in its arguments, the above expression will give different behavior compared to minimizing $KL(p||q)$ known as M-projection or moment projection:

$$(\text{M} - \text{projection}) : \quad q(x) = \arg \min_{q \in Q} KL(p||q) \quad (254)$$

Both the I-projection and the M-projection are projections of a probability distribution $p(x)$ onto a set of distributions Q . For I-projection, $q(x)$ will typically under-estimate the support of $p(x)$ and will lock onto one of its modes. This is due to $q(x) = 0$ whenever $p(x) = 0$ to make sure KL divergence stays finite. For M-projection, $q(x)$ will typically over-estimate the support

of $p(x)$ and will cover all of its modes. This is due to $q(x) > 0$ whenever $p(x) > 0$ to make sure KL divergence stays finite. The I-projection is useful in setting up information geometry because of the following inequality:

$$KL(q||p) \geq KL(q||p^*) + KL(p^*||p) \quad (255)$$

The inequality can be interpreted as information-geometric version of Pythagoras' triangle inequality theorem, where KL divergence is viewed as squared distance in Euclidean space.

If we are interested in measuring the distance between two multivariate Gaussian distributions with means μ_1 and μ_2 and covariances Σ_1 and Σ_2 , the KL divergence can be computed in closed form:

$$\begin{aligned} KL(p||q) &= \int p(x) \log \frac{p(x)}{q(x)} = \int [\log p(x) - \log q(x)] p(x) dx \\ &= \int \left[\frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) \right] p(x) dx \\ &= \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} - \frac{1}{2} \text{Tr}\{E[(x - \mu_1)(x - \mu_1)^T \Sigma_1^{-1}]\} + \frac{1}{2} E[(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)] \\ &= \frac{1}{2} \log \frac{|\Sigma_2|}{|\Sigma_1|} - \frac{1}{2} \text{Tr}\{I_d\} + \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) + \frac{1}{2} \text{Tr}\{\Sigma_2^{-1} \Sigma_1\} \\ &= \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - d + \text{Tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right] \end{aligned} \quad (256)$$

5.2.3 Mutual Information

Mutual information (MI) is a measure of overlap of random variables: amount of information one random variable contains about another random variable. Mutual information measures how similar the joint distribution $p(x, y)$ is compared to the factorized distribution $p(x)p(y)$ when the two variables are independent:

$$I(X; Y) = KL(p(x, y)||p(x)p(y)) = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \quad (257)$$

where $I(X; Y) \geq 0$ with equality iff the variables are independent: $p(x, y) = p(x)p(y)$. Mutual information between X and Y can be interpreted as the reduction in uncertainty about X after observing Y or by symmetry, the reduction in uncertainty about Y after observing X :

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = H(X) + H(Y) - H(X, Y) \quad (258)$$

The relationship between entropy and MI is captured in a Venn diagram in Figure 55. Note that entropy can be viewed as self-information $H(X) = I(X; X)$. The MI is the expected value of pointwise mutual information:

$$PMI(x, y) = \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)} \quad (259)$$

This can be interpreted as the amount we learn from updating the prior $p(y)$ into the posterior $p(y|x)$. In fact, we can express MI as follows:

$$I(X; Y) = \sum_x p(x) \sum_y p(y|x) \log \frac{p(y|x)}{p(y)} = E_{p(x)}[KL(p(y|x)||p(y))] \quad (260)$$

Therefore, the MI between X and Y is equivalent to the expected KL distance between the posterior distribution $p(y|x)$ and the prior $p(y)$.

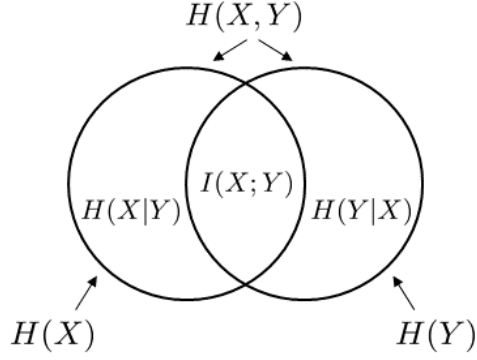


Figure 55: Relationship between entropy and mutual information

Theorem 5.4 (*Non-decreasing property of MI*). *As we add more random variables, the MI can increase or stay the same:*

$$I(X; Y_1) \leq I(X; Y_1, Y_2) \leq I(X; Y_1, \dots, Y_n) \quad (261)$$

where equality holds when Y_i 's are independent of X .

We can compute a variational lower bound on MI using the Gibbs inequality.

Theorem 5.5 (*Variational MI bound*). *Let $p(x)$ be the original posterior distribution and $q(x)$ be the variational approximation to the posterior. The, the following inequality holds:*

$$I(X; Y) = \max_q \left[H_p(x) + E_p[\log q(x)] \right] \quad (262)$$

Proof. By non-negativity of KL divergence, we can obtain the Gibbs inequality.

$$\text{KL}(p||q) = \sum_x p(x) \log p(x)q(x) = -H(p) + H(p, q) \geq 0 \quad (263)$$

As a result, we have

$$-H_p(x) \geq -H_{p,q}(x) \quad (264)$$

$$I(X; Y) = H(X) - H_p(X|Y) \geq H(X) - H_{p,q}(X) = H_p(x) + E_p[\log q(x)] \quad (265)$$

Theorem 5.6 (*Data Processing Inequality*). *Let X, Y, Z form the following Markov chain: $X \rightarrow Y \rightarrow Z$, then the following inequality holds:*

$$I(X; Y) \geq I(X; Z) \quad (266)$$

Proof. By the chain rule, we can expand mutual information in two different ways:

$$I(X; Y, Z) = I(X; Z) + I(X; Y|Z) = I(X; Y) + I(X; Z|Y) \quad (267)$$

Since X is conditionally independent of Z given Y , we have $I(X; Z|Y) = 0$ and therefore

$$I(X; Y) = I(X; Z) + I(X; Y|Z) \quad (268)$$

Since $I(X; Y|Z) \geq 0$, we have the inequality: $I(X; Y) \geq I(X; Z)$. We can use the data processing inequality to show that *sufficient statistics* preserves mutual information. In particular consider

the data generating process: $\theta \rightarrow X \rightarrow T(X)$. Given distribution parameters θ , we generate data by sampling from $f_\theta(x)$ and then we compute a statistic $T(X)$. The statistic $T(X)$ is called sufficient for θ if it contains all the information in X about θ , i.e. the data processing inequality is satisfied with equality:

$$I(\theta; X) = I(\theta; T(X)) \quad (269)$$

Hence sufficient statistics compresses the information about θ using sampled data.

For a bivariate Gaussian distribution we can compute $I(X; Y)$ as follows. Let $\rho = \text{cov}(X, Y)/\sigma_x\sigma_y$, and $\sigma^2 = \text{var}(X) = \text{var}(Y)$ then $h(X) = h(Y) = \frac{1}{2} \log(2\pi e)\sigma^2$ and

$$h(X, Y) = \frac{1}{2} \log [(2\pi e)^2 |\Sigma|] = \frac{1}{2} \log(2\pi e)^2 \sigma^4 (1 - \rho^2) \quad (270)$$

Therefore,

$$I(X; Y) = h(X) + h(Y) - h(X, Y) = -\frac{1}{2} \log(1 - \rho^2) \quad (271)$$

5.3 Imbalanced Learning

Most classification algorithms will only perform optimally when the number of samples in each class is roughly the same. Highly skewed datasets where the minority class is outnumbered by one or more classes commonly occur in fraud detection, medical diagnosis and computational biology. One way of addressing this issue is by re-sampling the dataset to offset the imbalance and arrive at a more robust and accurate decision boundary. Re-sampling techniques can be broadly divided into four categories: undersampling the majority class, over-sampling the minority class, combining over and under sampling, and creating ensemble of balanced datasets [17].

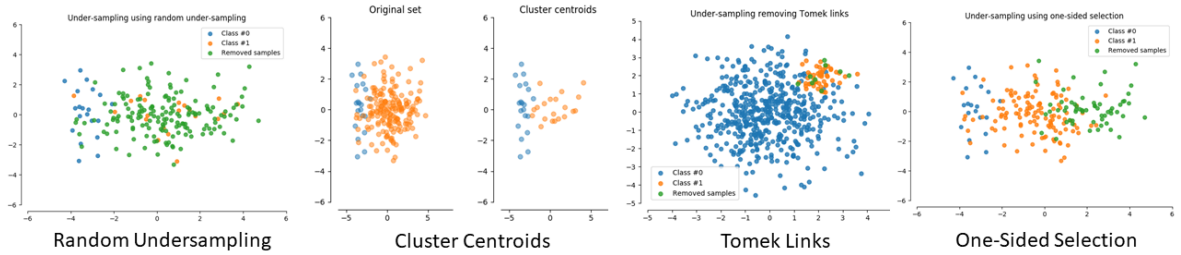


Figure 56: Four undersampling strategies: random, cluster centroids, Tomek links, one-sided selection.

Undersampling strategies. Undersampling methods remove data from the majority class of the original dataset as shown in Figure 56. Random Under Sampler simply removes data points from the majority class uniformly at random. Cluster Centroids is a method that replaces cluster of samples by the cluster centroid of a K-means algorithm, where the number of clusters is set by the level of undersampling. Tomek links remove unwanted overlap between classes where Tomek links are removed until all minimally distanced nearest neighbor pairs are of the same class. A Tomek link is defined as follows: given an instance pair (x_i, x_j) , where $x_i \in S_{min}$, $x_j \in S_{maj}$ and $d(x_i, x_j)$ is the distance between x_i and x_j , then the (x_i, x_j) pair is called a Tomek link if there's no instance x_k such that $d(x_i, x_k) < d(x_i, x_j)$ or $d(x_j, x_k) < d(x_i, x_j)$. In this way, if two instances form a Tomek link then either one of these instances is noise or both are near a border. Therefore one can use Tomek links to clean up overlap between classes. By removing overlapping examples, one can establish well-defined clusters in the training set and

lead to improved classification performance. The One Sided Selection (OSS) method selects a representative subset of the majority class E and combines it with the set of all minority examples S_{min} to form $N = \{E \cup S_{min}\}$. The reduced set N is further processed to remove all majority class Tomek links.

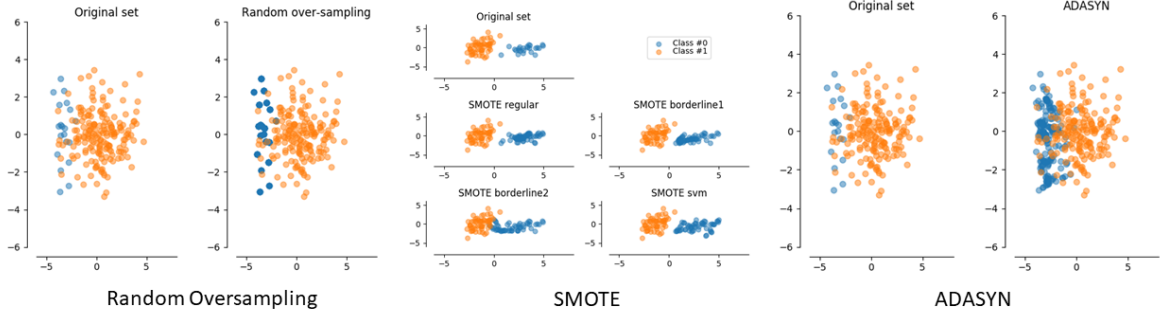


Figure 57: Three oversampling strategies: random, SMOTE, ADASYN.

Oversampling strategies. Oversampling methods append data to the minority class of the original dataset as shown in Figure 57. Random Over Sampler simply adds data points to the minority class uniformly at random. Synthetic Minority Oversampling Technique (SMOTE) generates synthetic examples by finding K nearest neighbors in the feature space and generating a new data point along the line segments joining any of the K minority class nearest neighbors. Synthetic samples are generated in the following way: take the difference between the feature vector (sample) under consideration and its nearest neighbor, multiply this difference by a random number between 0 and 1 and add it to the feature vector under consideration thus augmenting the dataset with a new data point. Adaptive Synthetic Sampling (ADASYN) uses a weighted distribution for different minority class examples according to their level of difficulty in learning, where more synthetic data is generated for minority class examples that are harder to learn. As a result, the ADASYN approach improves learning of imbalanced dataset in two ways: reducing the bias introduced by class imbalance and adaptively shifting the classification decision boundary toward the difficult examples.

It's possible to combine over-sampling and under-sampling techniques into a hybrid strategy. Common examples include SMOTE and Tomek links or SMOTE and Edited Nearest Neighbors (ENN). Additional ways of learning on imbalanced datasets include weighing training instances, introducing different misclassification costs for positive and negative examples and bootstrapping.

5.4 Ensemble Methods

Ensemble methods are meta-algorithms that combine several machine learning techniques into one predictive model in order to decrease variance (bagging), bias (boosting), or improve predictions (stacking). Ensemble methods can be divided into two groups: *sequential* ensemble methods where the base learners are generated sequentially (e.g. AdaBoost) and *parallel* ensemble methods where the base learners are generated in parallel (e.g. Random Forest). The basic motivation of sequential methods is to exploit the dependence between the base learners since the overall performance can be boosted by weighing previously mislabeled examples with higher weight. The basic motivation of parallel methods is to exploit independence between the

base learners since the error can be reduced dramatically by averaging.

Most ensemble methods use a single base learning algorithm to produce homogeneous base learners, i.e. learners of the same type leading to *homogeneous ensembles*. There are also some methods that use heterogeneous learners, i.e. learners of different types, leading to *heterogeneous ensembles*. In order for ensemble methods to be more accurate than any of its individual members the base learners have to be as accurate as possible and as diverse as possible.

5.4.1 Bagging

Bagging stands for bootstrap aggregation. One way to reduce the variance of an estimate is to average together multiple estimates. For example, we can train M different trees f_m on different subsets of the data (chosen randomly with replacement) and compute the ensemble:

$$f(x) = \frac{1}{M} \sum_{m=1}^M f_m(x) \quad (272)$$

Bagging uses bootstrap sampling to obtain the data subsets for training the base learners. For aggregating the outputs of base learners, bagging uses voting for classification and averaging for regression.

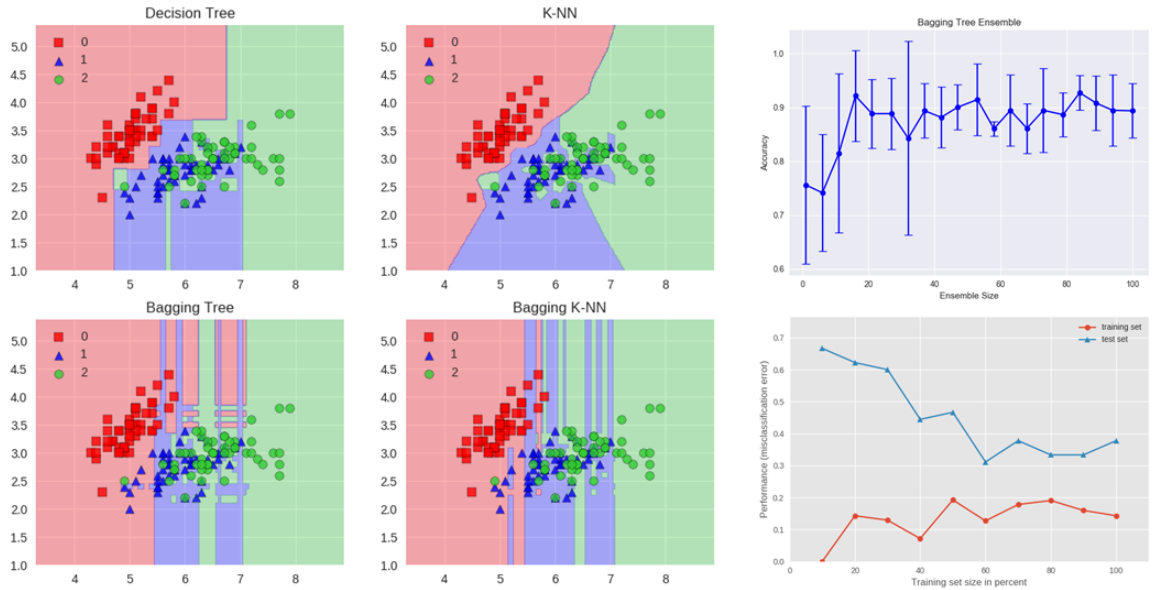


Figure 58: Bagging Ensemble with decision tree and k-NN as base estimator.

Figure 58 shows the decision boundary of a decision tree and k-NN classifiers along with their bagging ensembles applied to the Iris dataset. The decision tree shows axes parallel boundaries while the $k = 1$ nearest neighbors fits closely to the data points. The bagging ensembles were trained using 10 base estimators with 0.8 subsampling of training data and 0.8 subsampling of features. The decision tree bagging ensemble achieved higher accuracy in comparison to k-NN bagging ensemble because k-NN are less sensitive to perturbation on training samples and therefore they are called *stable learners*. Combining stable learners is less advantageous since the ensemble will not help improve generalization performance. The figure also shows how the test accuracy improves with the size of the ensemble. Based on cross-validation results, we can

see the accuracy increases until approximately 20 base estimators and then plateaus afterwards. Thus, adding base estimators beyond 20 only increases computational complexity without accuracy gains for the Iris dataset. The figure also shows learning curves for the bagging tree ensemble. We can see an average error of 0.15 on the training data and a U-shaped error curve for the testing data. The smallest gap between training and test errors occurs at around 80% of the training set size.

A commonly used class of ensemble algorithms are forests of randomized trees. In **random forests**, each tree in the ensemble is built from a sample drawn with replacement (i.e. a bootstrap sample) from the training set. In addition, instead of using all the features, a random subset of features is selected further randomizing the tree. As a result, the bias of the forest increases slightly but due to averaging of less correlated trees, its variance decreases resulting in an overall better model.

In **extremely randomized trees** algorithm randomness goes one step further: the splitting thresholds are randomized. Instead of looking for the most discriminative threshold, thresholds are drawn at random for each candidate feature and the best of these randomly-generated thresholds is picked as the splitting rule. This usually allows to reduce the variance of the model a bit more, at the expense of a slightly greater increase in bias.

5.4.2 Boosting

Boosting refers to a family of algorithms that are able to convert weak learners to strong learners. The main principle of boosting is to fit a sequence of weak learners (models that are only slightly better than random guessing, such as small decision trees) to weighted versions of the data, where more weight is given to examples that were mis-classified by earlier rounds. The predictions are then combined through a weighted majority vote (classification) or a weighted sum (regression) to produce the final prediction. The principal difference between boosting and the committee methods such as bagging is that base learners are trained in sequence on a weighted version of the data.

Algorithm 11 describes the most widely used form of boosting algorithm called **AdaBoost** which stands for adaptive boosting [11].

Algorithm 11 AdaBoost

- 1: Init data weights $\{w_n\}$ to $1/N$
 - 2: **for** $m = 1$ to M **do**
 - 3: fit a classifier $y_m(x)$ by minimizing weighted error function J_m :
 - 4: $J_m = \sum_{n=1}^N w_n^{(m)} 1[y_m(x_n) \neq t_n]$
 - 5: compute $\epsilon_m = \sum_{n=1}^N w_n^{(m)} 1[y_m(x_n) \neq t_n] / \sum_{n=1}^N w_n^{(m)}$
 - 6: evaluate $\alpha_m = \log\left(\frac{1-\epsilon_m}{\epsilon_m}\right)$
 - 7: update the data weights: $w_n^{(m+1)} = w_n^{(m)} \exp\{\alpha_m 1[y_m(x_n) \neq t_n]\}$
 - 8: **end for**
 - 9: Make predictions using the final model: $Y_M(x) = \text{sign}\left(\sum_{m=1}^M \alpha_m y_m(x)\right)$
-

We see that the first base classifier $y_1(x)$ is trained using weighting coefficients $w_n^{(1)}$ that are all equal. In subsequent boosting rounds, the weighting coefficients $w_n^{(m)}$ are increased for

data points that are misclassified and decreased for data points that are correctly classified. The quantity ϵ_m represents a weighted error rate of each of the base classifiers. Therefore, the weighting coefficients α_m give greater weight to the more accurate classifiers.

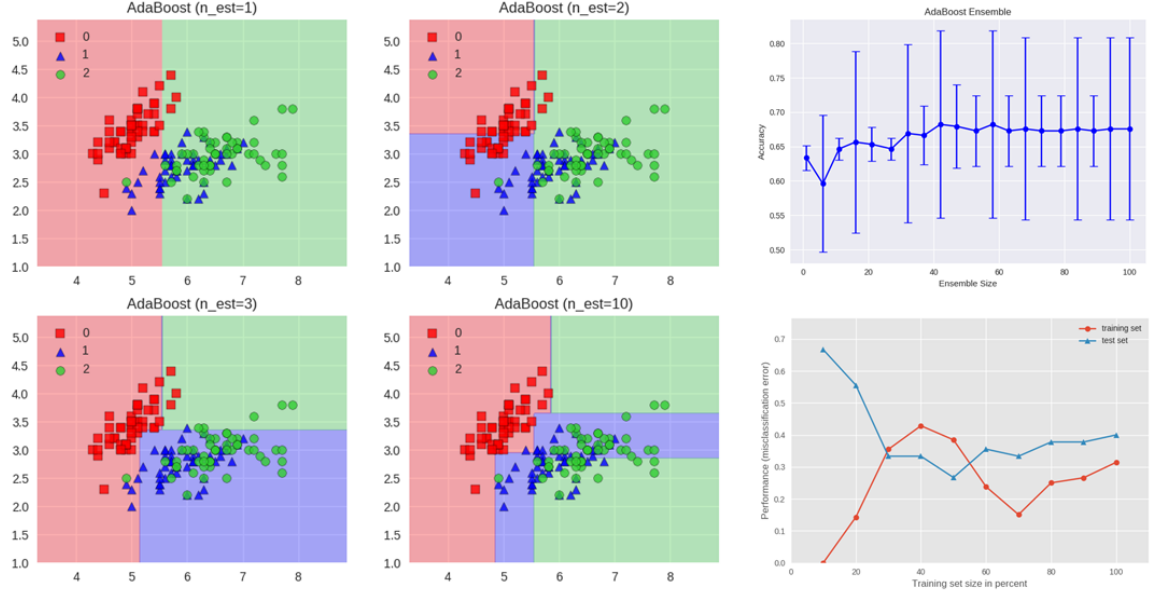


Figure 59: Boosting Ensemble with decision tree as base estimator.

The AdaBoost algorithm is illustrated in Figure 59. Each base learner consists of a decision tree with depth 1, thus classifying the data based on a feature threshold that partitions the space into two regions separated by a linear decision surface that is parallel to one of the axes. The figure also shows how the test accuracy improves with the size of the ensemble and the learning curves for training and testing data.

Gradient Tree Boosting is a generalization of boosting to arbitrary differentiable loss functions. It can be used for both regression and classification problems. Gradient Boosting builds the model in a sequential way:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x) \quad (273)$$

At each stage the decision tree $h_m(x)$ is chosen to minimize a loss function L given the current model $F_{m-1}(x)$:

$$F_m(x) = F_{m-1}(x) + \arg \min_h \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + h(x_i)) \quad (274)$$

Gradient Boosting attempts to solve this minimization problem numerically via steepest descent. The steepest descent direction is the negative gradient of the loss function evaluated at the current model F_{m-1} :

$$F_m(x) = F_{m-1}(x) + \gamma_m \sum_{i=1}^n \nabla_F L(y_i, F_{m-1}(x_i)) \quad (275)$$

where the step size γ_m is chosen using line search. The algorithms for regression and classification differ in the type of loss function used.

5.4.3 Stacking

Stacking is an ensemble learning technique that combines multiple classification or regression models via a meta-classifier or a meta-regressor. The base level models are trained based on complete training set then the meta-model is trained on the outputs of base level model as features. The base level often consists of different learning algorithms and therefore stacking ensembles are often heterogeneous. Algorithm 12 summarizes stacking.

Algorithm 12 Stacking

```

1: Input: training data  $D = \{x_i, y_i\}_{i=1}^m$ 
2: Output: ensemble classifier  $H$ 
3: Step 1: learn base-level classifiers
4: for  $t = 1$  to  $T$  do
5:   learn  $h_t$  based on  $D$ 
6: end for
7: Step 2: construct new data set of predictions
8: for  $i = 1$  to  $m$  do
9:    $D_h = \{x'_i, y_i\}$ , where  $x'_i = \{h_1(x_i), \dots, h_T(x_i)\}$ 
10: end for
11: Step 3: learn a meta-classifier
12: learn  $H$  based on  $D_h$ 
13: return  $H$ 

```

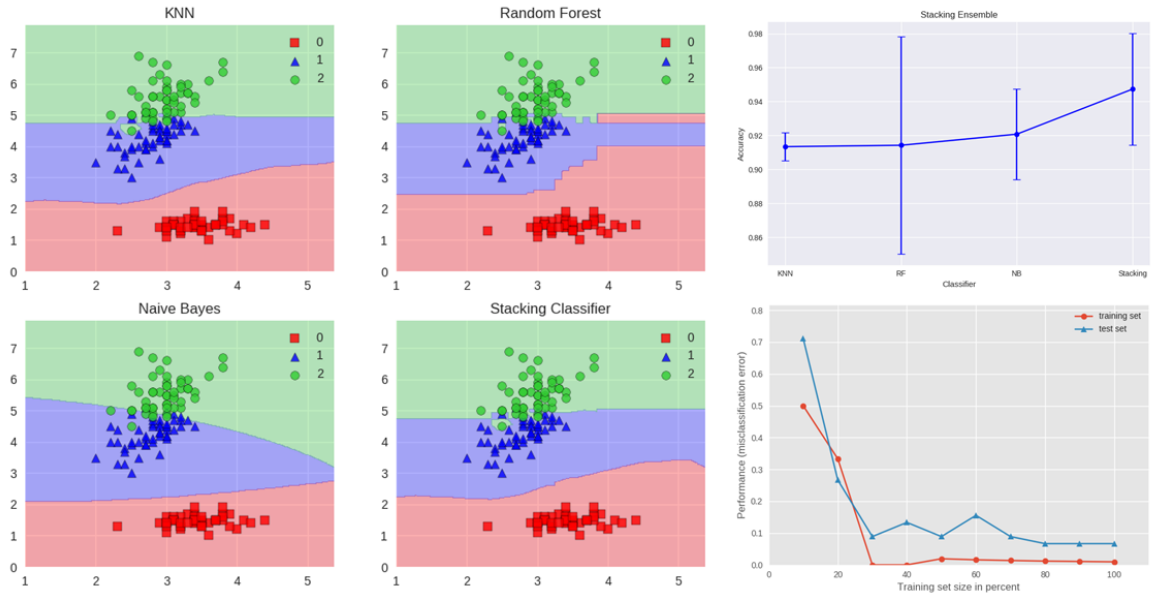


Figure 60: Stacking Ensemble with Logistic Regression as meta-classifier.

The stacking ensemble is illustrated in Figure 60. It consists of k-NN, Random Forest and Naive Bayes base classifiers whose predictions are combined by Logistic Regression as a meta-classifier. We can see the blending of decision boundaries achieved by the stacking classifier. The figure also shows that stacking achieves higher accuracy than individual classifiers and based on learning curves, it shows no signs of overfitting.

References

- [1] M. Becker, B. Hachey, B. Alex, and C. Grover. Optimising selective sampling for bootstrapping named entity recognition. In *ICML*, 2005.
- [2] M. Becker and M. Osborne. A two-stage method for active learning of statistical grammars. In *ACL*, 2002.
- [3] D. Blei and J. McAuliffe. Supervised topic models. In *NIPS*, pages 1–8, 2010.
- [4] D. Blei, A. Ng, and M. Jordan. Latent dirichlet allocation. In *JMLR*, pages 993–1022, 2003.
- [5] T. Campbell, M. Liu, B. Kulis, J. P. How, and L. Carin. Dynamic clustering via asymptotics of the dependent dirichlet process. In *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- [6] K. Chaloner and I. Verdinelli. Bayesian experimental design: a review. In *Statistical Science*, 1995.
- [7] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley, 2006.
- [8] N. de Freitas. Rao-blackwellized particle filtering for fault diagnosis. In *IEEE Aerospace Conference*, 2002.
- [9] F. Doshi, K. T. Miller, J. Van Gael, and Y. W. Teh. Variational inference for the Indian buffet process. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 12, 2009.
- [10] V. Federov. *Theory of Optimal Experiments*. Academic Press, 1972.
- [11] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. In *J. Comput. Syst. Sci.*, pages 119–139, 1997.
- [12] Y. Gal. Uncertainty in deep learning (phd thesis). In *arXiv*, 2016.
- [13] Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *ICML*, pages 1–12, 2016.
- [14] Y. Gal, R. Islam, and Z. Ghahramani. Deep bayesian active learning with image data. In *arXiv*, pages 1–8, 2017.
- [15] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. The MIT Press, 2016.
- [16] T. L. Griffiths and Z. Ghahramani. The indian buffet process: An introduction and review. In *J. Mach. Learn. Res.*, pages 1185–1224, 2011.
- [17] H. He and E. A. Garcia. Learning from imbalanced data. In *IEEE Trans. on Knowl. and Data Eng.*, pages 1263–1284, 2009.
- [18] M. Hoffman, F. R. Bach, and D. M. Blei. Online learning for latent dirichlet allocation. In *Advances in Neural Information Processing Systems*, pages 856–864, 2010.
- [19] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley. Stochastic variational inference. In *Journal of Machine Learning Research*, pages 1303–1347, 2013.
- [20] Kaggle. Sentiment analysis of movie reviews dataset. In *Kaggle Datasets*, 2015.
- [21] Koller and Friedman. *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press, 2009.
- [22] D. Lewis and W. Gale. A sequential algorithm for training text classifiers. In *ACM SIGIR*, 1994.
- [23] D. Lin, W. Grimson, and J. W. Fisher III. Construction of dependent dirichlet processes based on compound poisson processes. In *Neural Information Processing Systems*, 2010.

- [24] S. N. MacEachern. Dependent nonparametric processes. In *Proceedings of the Bayesian Statistical Science Section*, 1999.
- [25] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. In *arXiv*, 2013.
- [26] K. P. Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012.
- [27] G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher. An analysis of approximations for maximizing submodular set functions. In *Mathematical Programming*, pages 265–294, 1978.
- [28] F. Olsson. A literature of active machine learning in the context of natural language processing. In *Technical Report*, 2009.
- [29] B. Pang and L. Lee. A sentimental education: sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the ACL*, 2004.
- [30] A. Patil, D. Huard, and C. Fonnesbeck. PyMC: Bayesian stochastic modelling in python. In *Journal of Statistical Software*, pages 1–81, 2010.
- [31] J. Pennington, R. Socher, and C. Manning. Glove: Global vectors for word representation. In *EMNLP*, pages 1–12, 2014.
- [32] E. Ringger, P. McClanahan, R. Haertel, G. Busby, M. Carmen, J. Carroll, K. Seppi, and D. Lonsdale. Active learning for part-of-speech tagging: accelerating corpus annotation. In *ACL*, 2002.
- [33] T. Scheffer, C. Decomain, and S. Wrobel. Query by committee. In *IDA*, 2001.
- [34] J. Sethuraman. A constructive definition of Dirichlet priors. In *Statistica Sinica*, pages 639–650, 1994.
- [35] B. Settles. Active learning literature survey. In *Computer Sciences Technical Report*, pages 1–67, 2009.
- [36] H. S. Seung, M. Oppor, and H. Sompolinsky. Query by committee. In *ACM COLT*, 1992.
- [37] J. Snoek, H. Larochelle, and R. P. Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959, 2012.
- [38] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. In *JMLR*, pages 1929–1958, 2014.
- [39] Y. Teh, M. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet Processes. In *Journal of the American Statistical Association (JASA)*, 2005.
- [40] Y. W. Teh. Dirichlet processes. In *Encyclopedia of Machine Learning*. Springer, 2010.
- [41] S. Tong and D. Koller. Support vector machine active learning with applications to text classification. In *JMLR*, 2002.
- [42] J. Van Gael, Y. Saatci, Y. W. Teh, and Z. Ghahramani. Beam sampling for the infinite hidden Markov model. In *Proceedings of the International Conference on Machine Learning*, volume 25, 2008.
- [43] C. Wang, J. Paisley, and D. Blei. Online variational inference for the hierarchical dirichlet process. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 752–760, 2011.