

Supervised Learning: Find a mapping f from input features x to target outputs t using **labelled** training data (input-output pairs).

Unsupervised Learning: Find interesting patterns in the data, e.g., find clusters, where **no target labels are given** (only inputs are given).

Reinforcement Learning: Find a sequence of actions that produce a high reward signal in an environment.

Training Set: Used to learn the (*parameters* of) the mapping f from input features to the predicted output

Validation Set: Used to choose *hyperparameters* or model settings that are not provided by the learning algorithm

Test Set: Used to estimate how well the model will *generalize* to new data and how well the model will perform when deployed

Training:

k-Nearest Neighbours: No training necessary; just store the data

Decision Trees: Greddy method: find the optimum split at each node, where “optimum” is determined by a loss, e.g., information gain

Linear Regression and Classification: Find weights w that minimize a cost function, e.g., mean square error or average cross entropy loss across the training set

$x_j^{(i)}$ represents the j -th feature of the i -th training example

Hypothesis (model, predictor) is a function f that maps input x to output prediction y

Voroni diagram partitions space into regions closest to each data

kNN: k too small may overfit (sensitive to noise), k too large would underfit (fail to capture regularities)

$k < \sqrt{n}$, choose k with highest validation accuracy

We use test data to estimate generalization error

60% training, 20% validation and test

We normalize each dimension to zero mean and unit variance

Curse of dimensionality: data not dense in high dimensions, and hard to calculate distances. Most points are approximately equally distanced, squared Euclidean distance becomes dominated by the average behavior, making distances converge

Overfitting: High training accuracy, poor test accuracy; model too complex, small data set, noise in data; regularization, better hyperparameters, better model

Underfitting: Low training and test accuracy; model too simple, bad hyperparameters, bad model; adjust hyperparameters, regularization, normalization

Decision Tree Output: Most common value for discrete classification, or mean value for regression

Use greedy heuristic to construct decision tree based on (highest) information gain

More certain outcome has lower entropy, i.e. not uniformly distributed

Entropy characterizes uncertainty of random variable

Conditional Entropy characterizes uncertainty of random variable given another random variable

Information Gain measures informativeness of variable

Process: Pick a feature, compare all possible splits, choose the one with highest IG, repeat recursively

We may minimize squared error instead: $SE = \sum_{i=1}^N (f(x^{(i)}) - t^{(i)})^2$

A model is a set of assumptions (restrictions) about the structure of f

The model or architecture restricts the hypothesis f space

Goal of Linear Regression: find w (weights) and b (bias) minimizing $\mathcal{E}(w, b) = \mathcal{E}(\mathbf{w})$

Let X be matrix with $N \times (D + 1)$ shape with each $x^{(i)}$, then $\mathbf{y} = \mathbf{X}\mathbf{w}$ **Gradient Descent** is more general than matrix inversion for finding \mathbf{w}

Each GD update costs $O(ND)$, $\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} \mathcal{E}(\mathbf{w}) = \mathbf{w} - \frac{\alpha}{N} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{t})$, or $w_i \leftarrow w_i - \alpha \frac{\partial \mathcal{E}}{\partial w_i}$, where initially $\alpha = 0.002$, $w_i = b = 0$

If α too small, training is slow, if α large, diverges

We may generalize linear regression by mapping input feature to another space, by feature mapping (or basis expansion): $\psi(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^d$, and treat \mathbb{R}^d as input space, e.g. polynomial $\psi(x) = (1, x, \dots, x^M)^T$, then $y = \psi(x)^T \mathbf{w}$

Adjust the degree of polynomial to avoid underfitting or overfitting

regularizer or **penalty** is function that quali-

fies how much we prefer one hypothesis (penalize large weights)

Both polynomial degree M and λ are hyperparameters

λ large cause underfit, λ small cause overfit
We want each feature’s values to be indicators or orderable

ML Algorithm Design: 1. Choose model; 2. Define loss function; 3. Choose regularizer; 4. Fit model by minimizing regularized cost, using **optimization algorithm**.

$t = 1$ for positive example, $t = 0$ or -1 for negative example

For **Linear Binary Classification Model**, let $z = \mathbf{w}^T x$, then $y = 1$ if $z \geq 0$, else $y = 0$, decision boundary being hyperplane $z = 0$

Given small training set, go over each data point to find restrictions of w_i , if resion in weight space satisfying all constraints is non-empty, then this **feasible region** makes the problem (in)-**festible** and the training set **linearly separable**

Use square error loss function would cause correctly classified points with high loss; use indicator loss function would cause gradient be zero e.w.; use sigmoid plus square loss won’t update much for very wrongly classified example, i.e. poor gradient signal

We use sigmoid to smooth the output of each choice of \mathbf{w} between 0 and 1, then use cross entropy loss function to penalize the model for being wrong

$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) \mathbf{x}^{(i)}$, same as linear regression

hypothesis space \mathcal{H} set of all possible h a model can represent, i.e. all linear functions for linear regression

Members of \mathcal{H} with algorithm’s preference determine **inductive bias**

Algorithm is **parametric** if \mathcal{H} is defined by a set of parameters (linear and logistic regression, neutral network; counter-example kNN)

Previously we used **full batch gradient descent** of $\frac{\partial \mathcal{E}}{\partial \mathbf{w}}$, we may estimate $\frac{\partial \mathcal{E}}{\partial \mathbf{w}}$ using subset of data if N too large, via **stochastic gradient descent**: computing average of a small number of examples called **mini-batch** with a **batch size**.

inner loop is **iteration** (each mini-batch), outer is **epoch** (one update)

num.iter = num.epoch * $N/\text{batch_size}$

Accuracy: $\frac{1}{N} \sum_{i=1}^N \mathbb{I}[t^{(i)} = y^{(i)}]$

Error: $1 - \text{Accuracy}$

Euclidean:

$$\left\| x^{(a)} - x^{(b)} \right\|_2 = \sqrt{\sum_{j=1}^d (x_j^{(a)} - x_j^{(b)})^2}$$

Cosine Similarity:

$$\text{cosine}(x^{(a)}, x^{(b)}) = \frac{x^{(a)} \cdot x^{(b)}}{\left\| x^{(a)} \right\|_2 \left\| x^{(b)} \right\|_2}$$

Mean: $\bar{x} = \frac{1}{N} \sum_{i=1}^N x^{(i)} = \mu$

Variance: $\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \bar{x})^2$

Standard Deviation: $\sigma = \sqrt{\sigma^2}$

Normalization: $\hat{x}_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{\sigma_j}$

Entropy:

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

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

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

$$\begin{aligned} H(Y|X) &= \sum_{x \in X} p(x) H(Y|x) \\ &= -\sum_{\substack{x \in X \\ y \in Y}} p(x, y) \log_2 p(y|x) \end{aligned}$$

$$H(X, Y) = H(Y, X) = H(X|Y) + H(Y)$$

$$H \geq 0, H(Y|Y) = 0, H(Y|X) \leq H(Y)$$

$$H(Y|X) = H(Y) \text{ if } X \text{ and } Y \text{ ind.}$$

Information Gain:

$$IG(Y|X) = H(Y) - H(Y|X)$$

Linear Model:

$$y = f(x) = \sum_{j=1}^D w_j x_j + b = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

$$\text{where } \tilde{\mathbf{w}}_1 = b, \tilde{\mathbf{w}}_i = w_{i-1}, \tilde{\mathbf{x}}_1 = 1, \tilde{\mathbf{x}}_i = x_{i-1}$$

Squared error loss function:

$$\mathcal{L}(y, t) = \frac{1}{2} (y - t)^2$$

$y - t$ is called residual

Cost (average loss) function:

$$\mathcal{E}(w, b) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, t^{(i)})$$

$$= \frac{1}{2N} (\mathbf{y} - \mathbf{t})^T (\mathbf{y} - \mathbf{t})$$

$$= \frac{1}{2N} \|\mathbf{X}\mathbf{w} - \mathbf{t}\|_2^2$$

$$\frac{\partial \mathcal{E}}{\partial w} = \frac{1}{N} \sum_{j=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)}$$

$$\frac{\partial \mathcal{E}}{\partial b} = \frac{1}{N} \sum_{j=1}^N (y^{(i)} - t^{(i)})$$

$$\frac{\partial \mathcal{L}}{\partial w_j} = (y - t) x_j$$

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) &= (y - t) \mathbf{x} \\ \nabla_{\mathbf{w}} \mathcal{E}(\mathbf{w}) &= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) \mathbf{x}^{(i)} \\ &= \frac{1}{N} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{t}) \end{aligned}$$

$$\nabla_{\mathbf{w}} \mathcal{E}(\mathbf{w}) = 0: \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

L^2 Regularizer:

$$\mathcal{R}(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2 = \frac{1}{2} \sum_{j=1}^D w_j^2$$

Regularized cost function:

$$\mathcal{E}_{\text{reg}}(\mathbf{w}) = \mathcal{E}(\mathbf{w}) + \lambda \mathcal{R}(\mathbf{w})$$

Logistic / Sigmoid function:

$$y = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Both sigmoid and square loss:

$$\frac{\partial \mathcal{L}}{\partial w_j} = (y - t) y (1 - y) x_j$$

Cross Entropy Loss:

$$\begin{aligned} \mathcal{L}_{\text{CE}}(y, t) &= \begin{cases} -\log(y) & \text{if } t = 1 \\ -\log(1 - y) & \text{if } t = 0 \end{cases} \\ &= -t \log(y) - (1 - t) \log(1 - y) \end{aligned}$$

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial w_j} = (y - t) x_j$$

3rd column cont., batch size too large = expensive, too small = noisy
center and normalize to avoid ravines

Graph of kNN: Draw tangent line at each midpoint, connect regions with same class

Graph of Decision Tree: Draw axis-aligned lines to split regions

IG of Decision Tree: First calculate the entropy $H(Y)$ with probability (percentages) of each class, similarly for $H(Y|\text{right})$, $H(Y|\text{left})$, then $IG = H(Y) - p(\text{right})H(Y|\text{right}) - p(\text{left})H(Y|\text{left})$
Shape: # of rows \times # of columns

Linear Regression: Square error loss function, with average loss cost function

Logistic Regression: Sigmoid function, with cross entropy loss function

Decision Tree: Greedy method, with information gain loss function

kNN: No training, with Euclidean distance for lower dimension D

Multi-class Classification

One-hot vector $t \in \mathbb{R}^K$ has 1 at the index of the correct class, 0 elsewhere, for K targets

For $y \in \mathbb{R}^K$, $\sum_{k=1}^K y_k = 1$

For input $x \in \mathbb{R}^D$, do linear prediction for each z_i , $z = Wx + b$, then $y_k = \text{softmax}(z)_k = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$, their sum is 1, each $y_k \in [0, 1]$

$$z = Wx + b$$

$$y = \text{softmax}(z) = \frac{e^z}{\sum_{j=1}^K e^{z_j}}$$

$$\mathcal{L}_{\text{CE}}(y, t) = -\sum_{k=1}^K t_k \log(y_k) = -t^T \log(y)$$

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial w_k} = \frac{\partial \mathcal{L}_{\text{CE}}}{\partial z_k} \frac{\partial z_k}{\partial w_k} = (y_k - t_k)x$$

$$\mathcal{E}(W, b) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{CE}}(y^{(i)}, t^{(i)})$$

$$w_k \leftarrow w_k - \frac{\alpha}{N} \sum_{i=1}^N (y_k^{(i)} - t_k^{(i)})x^{(i)}$$

Linear models can only separate data into 2 half-planes, does not work with XOR

Neural Network

Multi-Class: x_1, \dots, x_D as inputs, w_{ij} connecting x_j to y_i as connection

b_1, \dots, b_K as bias, $f = \text{softmax}$, $y = \text{softmax}(Wx + b)$ as output

Multi-Layer: Input layer, hidden layers, output layer (input layer doesn't count in n -layer n.n.)

Input size: num of features, hidden size: hyper-par., output size: num of classes

One layer:

$$h_j = f(\sum_{i=1}^D w_{j,i}^{(1)} x_i + b_j^{(1)})$$

$$h = f(W^{(1)}x + b^{(1)})$$

$$z_k = \sum_{j=1}^K w_{k,j}^{(2)} h_j + b_k^{(2)}$$

$$z = W^{(2)}h + b^{(2)}$$

$$y = \text{softmax}(z)$$

Multiple layer:

$$h^{(1)} = f(W^{(1)}x + b^{(1)})$$

$$h^{(\ell)} = f(W^{(\ell)}h^{(\ell-1)} + b^{(\ell)})$$

$$z = W^{(L)}h^{(L-1)} + b^{(L)}$$

$$y = \text{softmax}(z)$$

Common f : Sigmoid, Tanh, ReLU

Sigmoid: problematic due to gradient signal at extreme inputs, pos only

Tanh: Sigmoid but centered at 0, pos and neg
ReLU: $f(x) = \max(0, x)$, often used, pos, problematic if bias too large and neg (activation always 0)

N.N.: learning features s.t. becomes lin. sep. after $L - 1$ layers; final layer as linear classifier
Expressive power: Deep linear networks with

no activ. func. have same expr. power as linear regression

is universal approximator with nonlinear activ. func., e.g. single layer with 2^D hidden units

Backpropagation

For $\mathcal{L} = \frac{1}{2}(\sigma(wx + b) - t)^2$,

$$\frac{\partial \mathcal{L}}{\partial w} = (\sigma(wx + b) - t)\sigma'(wx + b)x$$

$$\frac{\partial \mathcal{L}}{\partial b} = (\sigma(wx + b) - t)\sigma'(wx + b)$$

So, more efficiently $\frac{\partial \mathcal{L}}{\partial y} = y - t$, $\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial y} \sigma'(z)$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial z} x, \frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z}$$

Let \bar{y} denote $\frac{d\mathcal{L}}{dy}$ called error signal

Computing predictions: $z = wx + b$, $y = \sigma(z)$, $\mathcal{L} = \frac{1}{2}(y - t)^2$

Computing gradients: $\bar{y} = y - t$, $\bar{z} = \bar{y}\sigma'(z)$, $\bar{w} = \bar{z}x$, $\bar{b} = \bar{z}$

Multiclass Logistic Regression with 2 features and 2 classes (computation graph):

$x_1, x_2 \rightarrow z_1, z_2 \rightarrow y_1, y_2 \rightarrow \mathcal{L}$,

$w_{11}, b_1, w_{12} \rightarrow z_1; w_{21}, b_2, w_{22} \rightarrow z_2; t_1, t_2 \rightarrow \mathcal{L}$

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

$$\bar{t} = \bar{x} \frac{dx}{dt} + \bar{y} \frac{dy}{dt}$$

$$\text{E.g. } \frac{\partial \mathcal{L}}{\partial b_1^{(1)}} = \left(\frac{\partial \mathcal{L}}{\partial y_1} \frac{\partial y_1}{\partial h_1} + \frac{\partial \mathcal{L}}{\partial y_2} \frac{\partial y_2}{\partial h_1} \right) \frac{dh_1}{dz_1} \frac{\partial z_1}{\partial b_1^{(1)}}$$

Forward pass: $z = W^{(1)}x + b^{(1)}$, $h = \sigma(z)$, $y = W^{(2)}h + b^{(2)}$, $\mathcal{L} = \frac{1}{2} \sum_k (y_k - t_k)^2 = \frac{1}{2} ||y - t||^2$

Backward pass: $\bar{\mathcal{L}} = 1$, $\bar{y} = \bar{\mathcal{L}}(y - t)$, $\bar{W}^{(2)} = \bar{y}h^T$, $\bar{b}^{(2)} = \bar{y}$, $\bar{h} = W^{(2)T}\bar{y}$, $\bar{z} = \bar{h} \circ \sigma'(z)$

If $\bar{z}_j = \sum_k \bar{y}_k \frac{\partial y_k}{\partial z_j}$, then $\bar{z} = \frac{\partial y}{\partial z}^T \bar{y}$, with

Jacobian matrix $\frac{\partial y}{\partial z}$

$$\text{For softmax, } \frac{\partial y_i}{\partial z_j} = \begin{cases} y_i(1 - y_j) & \text{if } i = j \\ -y_i y_j & \text{if } i \neq j \end{cases}$$

Bias Variance Decomposition
Decompose generalization error into: variance, bias, and irreducible error

Variance: Error from sensitivity to small fluctuations in the training data. Low Variance: less spread out. High Variance = overfitting

Bias: Error from poor assumptions in the model. Low Bias: more towards the center. High Bias = underfitting

Irreducible Error (Bayes Error): Error due to noise in the problem

Treat hypothesis y as a random variable. Lemma. Best prediction is $y_* = \mathbb{E}[t|x]$.

Decompose $\mathbb{E}[(y - t)^2|x]$ ignoring $|x$:
 $\mathbb{E}[(y - t)^2] = \mathbb{E}[(y - y_*)^2] + \text{Var}(t)$

$$= \mathbb{E}[y_* - y]^2 + \text{Var}(y - y_*) + \text{Var}(t)$$

I.e., expected prediction error = bias (perform on average), variance (amount of variability in predictions), and Bayes error (irreducible error

from data)

Ensembling methods: combine multiple models performs better than individual models

begging (bootstrap aggregation): train independently on different subsets of data, then average the predictions

Bayes error: unchanged; Bias: unchanged; Variance: reduced by factor $1/m$

Take \mathcal{D} with n examples; Generate m new datasets, each sample n training examples from \mathcal{D} with replacement; Averaging the predictions

$y_{\text{bagged}} = \mathbb{I}(\sum_{i=1}^m \frac{y_i}{m} > 0.5)$

Datasets not indep., not precisely $1/m$ reduc. Random forests = bagged decision trees, but choose a random set of d input features for each node of decision tree

Decreases correlation between classifiers by adding randomness

Probabilistic Classifiers
For binary classification, let Y be the hypothesis, $Y \sim \text{Bernoulli}(\theta)$ for unknown $\theta \in [0, 1]$

$p(y_i|\theta) = \theta^{y_i}(1 - \theta)^{1 - y_i}$

If y_i are i.i.d. Bernoullis, then $p(y_1, \dots, y_N|\theta) = \prod_{i=1}^N p(y_i|\theta) = \prod_{i=1}^N \theta^{y_i}(1 - \theta)^{1 - y_i}$

We may find θ by given y_i and maximizing Likelihood function $L(\theta) = p(y_1, \dots, y_N|\theta)$

maximum likelihood criterion:
pick $\hat{\theta}_{ML} = \arg \max_{\theta \in [0, 1]} L(\theta)$

Work with log-likelihoods in practice: $\ell(\theta) = \log L(\theta) = \sum_{i=1}^N y_i \log \theta + (1 - y_i) \log(1 - \theta)$

Let $N_H = \sum_i y_i$, $N_T = N - \sum_i y_i$, then set $\frac{d\ell}{d\theta} = \frac{N_H}{\theta} - \frac{N_T}{1 - \theta} = 0$

Obtain $\hat{\theta}_{ML} = \frac{N_H}{N_H + N_T} = \frac{N_H}{N}$

maximum likelihood estimation: define a model that assigns a probability to a dataset, maximize the likelihood (minimize cross-entropy)

$\log p(y_i, x_i; \theta) = (1 - x_i)(y_i \log \theta_0 + (1 - y_i) \log(1 - \theta_0)) + x_i(y_i \log \theta_1 + (1 - y_i) \log(1 - \theta_1)) + (x_i \log \pi + (1 - x_i) \log(1 - \pi))$

$\log p(\mathcal{D}; \theta) = \sum_{i=1}^N \log p(y_i, x_i; \theta)$

Decompose summation into 3 parts, take derivative of each part, set derivative to zero, solve for θ_0, θ_1, π

$$\begin{aligned} \text{E.g.: } p(\mathcal{D}|\theta) &= \prod_{i=1}^N p(x_1^{(i)}, x_2^{(i)}, y^{(i)}; \theta) \\ &= \prod_{i=1}^N p(y^{(i)}|x_1^{(i)}, x_2^{(i)})p(x_1^{(i)}, x_2^{(i)}) \\ &= \prod_{i=1}^N (\theta_{x_1^{(i)}, x_2^{(i)}})^{y^{(i)}} (1 - \theta_{x_1^{(i)}, x_2^{(i)}})^{1 - y^{(i)}} \end{aligned}$$

Then take log, decompose into 3 sums, for $\theta_{0,0}$, consider the first 2 sums only, set derivative to zero,

$$\frac{\partial \ell}{\partial \theta_{0,0}} = \frac{\partial}{\partial \theta_{0,0}} \sum \mathbb{I}[x_1^{(i)} = 0, x_2^{(i)} = 0, y^{(i)} = 1] \log \theta_{0,0} + \mathbb{I}[x_1^{(i)} = 0, x_2^{(i)} = 0, y^{(i)} = 0] \log(1 - \theta_{0,0})$$

$$= \frac{N_{0,0;p}}{\theta_{0,0}} - \frac{N_{0,0;N}}{1 - \theta_{0,0}}$$

$$\theta_{0,0}^* = \frac{N_{0,0;p}}{N_{0,0}}$$

$$\ell(\pi) = N_{0,0} \log \pi_{0,0} + \dots + N_{1,1} \log(1 - \dots - \pi_{1,0}) + C_0$$

$$\frac{\partial \ell}{\partial \pi_{0,0}} = 0 = \frac{N_{0,0}}{\pi_{0,0}} - \frac{N_{1,1}}{1 - \pi_{0,0} - \pi_{0,1} - \pi_{1,0}}$$

$$N_{0,0}\pi_{1,1} = N_{1,1}\pi_{0,0}, N_{0,1}\pi_{1,1} = N_{1,1}\pi_{0,1}$$

$$N_{1,0}\pi_{1,1} = N_{1,1}\pi_{1,0}$$

Discriminative classifier: learns a mapping from inputs to outputs, e.g. logistic regression, neural networks; model $p(c|x)$ directly (estimate parameters directly from labelled examples)

Generative model: Model $p(x, c), p(x|c)$, i.e. distribution of inputs characteristic of class

(Bayes classifier)
Bayes rule: $p(y|x) = \frac{p(y)}{p(x)} p(x|y)$

E.g. sample c from Bernoulli $p(c)$; sample x_1, \dots, x_D from $p(x|c)$ ($c = 1, c = 0$)

Learn via MLE
Discriminative: Given x , predict $p(c|x)$

Generative: Model $p(x, c)$, given $p(c)$ can compute $p(x|c)$ Binary Bag-of-words Features: $x_i = 1$ if word i appears in the document, 0 otherwise

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

Posterior = $\frac{\text{Class likelihood} \times \text{Class Prior}}{\text{Evidence}}$

If we want to compare $p(c = 0|x)$ with $p(c = 1|x)$, it suffices to compare $p(x|c = 0)p(c = 0)$ with $p(x|c = 1)p(c = 1)$

Evidence:
 $p(x) = \sum_c p(x|c)p(c)$

$= p(x|c = 1)p(c = 1) + p(x|c = 0)p(c = 0)$

$p(c, x_1, \dots, x_D)$ is enough to obtain $p(c)$ and $p(x|c)$ (using $2^{D+1} - 1$ entries)

Naive Bayes classifier assumes x_i are conditionally independent given c

$p(x_1, \dots, x_D|c) = \prod_{i=1}^D p(x_i|c)$

$p(c, x_1, \dots, x_D) = p(c) \prod_{i=1}^D p(x_i|c)$

$P(c = 1) = \pi$, $P(x_j = 1|c = i) = \theta_{ji}$

A directed graphical model (Bayesian network): joint distr. factorizes as a product of condi. distr. for each variable given parent

$\ell(\theta) = \sum_{i=1}^N \log p(c^{(i)}, x^{(i)})$

$= \sum_{i=1}^N \log(p(x^i|c^{(i)})p(c^{(i)}))$

$= \sum_{i=1}^N \log(p(c^{(i)}) + \prod_{j=1}^D p(x_j^{(i)}|c^{(i)}))$

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

where the first sum is Bernoulli log-likelihood of labels ($p(c) = \pi^c(1 - \pi)^{1-c}$, $\hat{\pi} = \text{pos/tot}$, $p(c) \text{ Ber.}$), second inner of feature x_j (decompose to learn each θ_{jc} separately, $p(x_j|c) \text{ Bern.}$, $\theta_{jc} = \text{num of word } j \text{ in pos/pos num}$)

$\theta_{jc} = p(x_j^{(i)} = 1|c) = \frac{x_j^{(i)}}{\theta_{jc}^{(i)}} (1 - \theta_{jc})^{1 - x_j^{(i)}}$

$\sum_{i=1}^N \log p(x_j^{(i)}|c^{(i)}) = \sum_{i=1}^N c^{(i)}$

$\{x_j^{(i)} \log \theta_{j1} + (1 - x_j^{(i)}) \log(1 - \theta_{j1})\}$

$+ \sum_{i=1}^N (1 - c^{(i)})$

$\{x_j^{(i)} \log \theta_{j0} + (1 - x_j^{(i)}) \log(1 - \theta_{j0})\}$

Inference: Compute numerator of $p(c|x) = \frac{p(c)p(x|c)}{\sum_{c'} p(c')p(x|c')}$

data sparsity: overfit when data too little
MLE: observations are R.V., parameters fixed

Bayesian: parameters as R.V., prior distr. $p(\theta)$, likelihood $p(\mathcal{D}|\theta)$

Update Posterior distr of θ : $p(\theta|\mathcal{D}) = \frac{p(\theta)p(\mathcal{D}|\theta)}{\int p(\theta')p(\mathcal{D}|\theta') d\theta'}$

E.g. coin $L(\theta) = p(\mathcal{D}|\theta) = \theta^{N_H}(1 - \theta)^{N_T}$, choose prior beta $p(\theta; a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1}(1 - \theta)^{b-1}$ (ignore normalization)

$p(\theta|\mathcal{D}) \propto [\theta^{a-1}(1 - \theta)^{b-1}] [\theta^{N_H}(1 - \theta)^{N_T}]$

which is beata with $a' = (a+)N_H, b' = (b+)N_T$

Posterior: $\mathbb{E}[\theta|\mathcal{D}] = \frac{N_H + a}{N_H + N_T + a + b}$

Maximum A-Posteriori Estimation: find most likely para. setting under posterior; converts par esti. to maxi. prob.

$\hat{\theta} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \max_{\theta} \log p(\theta) + \log p(\mathcal{D}|\theta) = \frac{N_H + a - 1}{N_H + N_T + a + b - 2}$