

Bayesian Methods & Multiclass

Mengye Ren

NYU

Oct 31, 2023

Announcement

- Schedule your project consultation soon.
- Use the provided template! (if your final report fails to use template then there will be marks off)
- Homework 3 is released and due Nov 14 11:59AM.

Recap

- Bayesian modeling adds a prior on the parameters.
- Models the distribution of parameters
- Bayes Rule:

$$p(y | x) = \frac{p(x | y)p(y)}{p(x)}$$

-

$$p(\theta | \mathcal{D}) = \frac{p(\mathcal{D} | \theta)p(\theta)}{p(\mathcal{D})}.$$

-

$$\underbrace{p(\theta | \mathcal{D})}_{\text{posterior}} \propto \underbrace{p(\mathcal{D} | \theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}}.$$

- Conjugate prior: Having the same form of distribution as the posterior.

Bayesian Point Estimates

- We have the posterior distribution $\theta | \mathcal{D}$.
- What if someone asks us to choose a single $\hat{\theta}$ (i.e. a point estimate of θ)?
- Common options:
 - **posterior mean** $\hat{\theta} = \mathbb{E}[\theta | \mathcal{D}]$
 - **maximum a posteriori (MAP) estimate** $\hat{\theta} = \arg \max_{\theta} p(\theta | \mathcal{D})$
 - Note: this is the **mode** of the posterior distribution

What else can we do with a posterior?

- Look at it: display uncertainty estimates to our client
- Extract a **credible set** for θ (a Bayesian confidence interval).
 - e.g. Interval $[a, b]$ is a 95% **credible set** if

$$\mathbb{P}(\theta \in [a, b] \mid \mathcal{D}) \geq 0.95$$

- Select a point estimate using **Bayesian decision theory**:
 - Choose a loss function.
 - Find action **minimizing expected risk w.r.t. posterior**

Bayesian Decision Theory

Bayesian Decision Theory

- Ingredients:
 - **Parameter space** Θ .
 - **Prior**: Distribution $p(\theta)$ on Θ .
 - **Action space** \mathcal{A} .
 - **Loss function**: $\ell : \mathcal{A} \times \Theta \rightarrow \mathbb{R}$.
- The **posterior risk** of an action $a \in \mathcal{A}$ is

$$\begin{aligned} r(a) &:= \mathbb{E}[\ell(\theta, a) \mid \mathcal{D}] \\ &= \int \ell(\theta, a) p(\theta \mid \mathcal{D}) d\theta. \end{aligned}$$

- It's the **expected loss under the posterior**.
- A **Bayes action** a^* is an action that minimizes posterior risk:

$$r(a^*) = \min_{a \in \mathcal{A}} r(a)$$

Bayesian Point Estimation

- General Setup:
 - Data \mathcal{D} generated by $p(y | \theta)$, for unknown $\theta \in \Theta$.
 - We want to produce a **point estimate** for θ .
- Choose:
 - **Prior** $p(\theta)$ on $\Theta = \mathbb{R}$.
 - **Loss** $\ell(\hat{\theta}, \theta)$
- Find **action** $\hat{\theta} \in \Theta$ that minimizes the **posterior risk**:

$$\begin{aligned}r(\hat{\theta}) &= \mathbb{E}[\ell(\hat{\theta}, \theta) | \mathcal{D}] \\ &= \int \ell(\hat{\theta}, \theta) p(\theta | \mathcal{D}) d\theta\end{aligned}$$

Important Cases

- Squared Loss : $\ell(\hat{\theta}, \theta) = (\theta - \hat{\theta})^2 \Rightarrow$ posterior mean
- Zero-one Loss: $\ell(\theta, \hat{\theta}) = \mathbb{1}[\theta \neq \hat{\theta}] \Rightarrow$ posterior mode
- Absolute Loss : $\ell(\hat{\theta}, \theta) = |\theta - \hat{\theta}| \Rightarrow$ posterior median
- Optimal decision depends on the loss function and the posterior distribution.
- Example: I have a card drawing from a deck of 2,3,3,4,4,5,5,5, and you guess the value of my card.
- mean: 3.875; mode: 5; median: 4

Bayesian Point Estimation: Square Loss

- Find **action** $\hat{\theta} \in \Theta$ that minimizes **posterior risk**

$$r(\hat{\theta}) = \int (\theta - \hat{\theta})^2 p(\theta | \mathcal{D}) d\theta.$$

- Differentiate:

$$\begin{aligned} \frac{dr(\hat{\theta})}{d\hat{\theta}} &= - \int 2(\theta - \hat{\theta}) p(\theta | \mathcal{D}) d\theta \\ &= -2 \int \theta p(\theta | \mathcal{D}) d\theta + 2\hat{\theta} \underbrace{\int p(\theta | \mathcal{D}) d\theta}_{=1} \\ &= -2 \int \theta p(\theta | \mathcal{D}) d\theta + 2\hat{\theta} \end{aligned}$$

Bayesian Point Estimation: Square Loss

- Derivative of posterior risk is

$$\frac{dr(\hat{\theta})}{d\hat{\theta}} = -2 \int \theta p(\theta | \mathcal{D}) d\theta + 2\hat{\theta}.$$

- First order condition $\frac{dr(\hat{\theta})}{d\hat{\theta}} = 0$ gives

$$\begin{aligned}\hat{\theta} &= \int \theta p(\theta | \mathcal{D}) d\theta \\ &= \mathbb{E}[\theta | \mathcal{D}]\end{aligned}$$

- The **Bayes action** for **square loss** is the posterior mean.

Interim summary

Recap and Interpretation

- The prior represents belief about θ before observing data \mathcal{D} .
- The posterior represents **rationally updated beliefs** after seeing \mathcal{D} .
- All inferences and action-taking are based on the posterior distribution.
- In the Bayesian approach,
 - No issue of justifying an estimator.
 - Only choices are
 - **family of distributions**, indexed by Θ , and
 - **prior distribution** on Θ
 - For decision making, we need a **loss function**.

Recap: Conditional Probability Models

Conditional Probability Modeling

- **Input space** \mathcal{X}
- **Outcome space** \mathcal{Y}
- **Action space** $\mathcal{A} = \{p(y) \mid p \text{ is a probability distribution on } \mathcal{Y}\}$.
- **Hypothesis space** \mathcal{F} contains prediction functions $f : \mathcal{X} \rightarrow \mathcal{A}$.
- **Prediction function** $f \in \mathcal{F}$ takes input $x \in \mathcal{X}$ and produces a **distribution** on \mathcal{Y}
- A **parametric family of conditional densities** is a set

$$\{p(y \mid x, \theta) : \theta \in \Theta\},$$

- where $p(y \mid x, \theta)$ is a density on **outcome space** \mathcal{Y} for each x in **input space** \mathcal{X} , and
- θ is a **parameter** in a [finite dimensional] **parameter space** Θ .
- This is the common starting point for either classical or Bayesian regression.

Classical treatment: Likelihood Function

- **Data:** $\mathcal{D} = (y_1, \dots, y_n)$
- The probability density for our data \mathcal{D} is

$$p(\mathcal{D} | x_1, \dots, x_n, \theta) = \prod_{i=1}^n p(y_i | x_i, \theta).$$

- For fixed \mathcal{D} , the function $\theta \mapsto p(\mathcal{D} | x, \theta)$ is the **likelihood function**:

$$L_{\mathcal{D}}(\theta) = p(\mathcal{D} | x, \theta),$$

where $x = (x_1, \dots, x_n)$.

- The **maximum likelihood estimator (MLE)** for θ in the family $\{p(y | x, \theta) | \theta \in \Theta\}$ is

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta \in \Theta} L_{\mathcal{D}}(\theta).$$

- MLE corresponds to ERM, if we set the loss to be the negative log-likelihood.
- The corresponding prediction function is

$$\hat{f}(x) = p(y | x, \hat{\theta}_{\text{MLE}}).$$

Bayesian Conditional Probability Models

Bayesian Conditional Models

- Input space $\mathcal{X} = \mathbb{R}^d$ Outcome space $\mathcal{Y} = \mathbb{R}$
- The Bayesian conditional model has two components:
 - A **parametric family of conditional densities**:

$$\{p(y | x, \theta) : \theta \in \Theta\}$$

- A **prior distribution** $p(\theta)$ on $\theta \in \Theta$.

The Posterior Distribution

- The **prior distribution** $p(\theta)$ represents our beliefs about θ before seeing \mathcal{D} .
- The **posterior distribution** for θ is

$$\begin{aligned} p(\theta | \mathcal{D}, x) &\propto p(\mathcal{D} | \theta, x) p(\theta) \\ &= \underbrace{L_{\mathcal{D}}(\theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}} \end{aligned}$$

- Posterior represents the **rationally updated beliefs** after seeing \mathcal{D} .
- Each θ corresponds to a prediction function,
 - i.e. the conditional distribution function $p(y | x, \theta)$.

Point Estimates of Parameter

- What if we want point estimates of θ ?
- We can use **Bayesian decision theory** to derive point estimates.
- We may want to use
 - $\hat{\theta} = \mathbb{E}[\theta | \mathcal{D}, x]$ (the posterior mean estimate)
 - $\hat{\theta} = \text{median}[\theta | \mathcal{D}, x]$
 - $\hat{\theta} = \arg \max_{\theta \in \Theta} p(\theta | \mathcal{D}, x)$ (the MAP estimate)
- depending on our loss function.

Back to the basic question - Bayesian Prediction Function

- Find a function takes input $x \in \mathcal{X}$ and produces a **distribution** on \mathcal{Y}
- In the frequentist approach:
 - Choose family of conditional probability densities (hypothesis space).
 - Select one conditional probability from family, e.g. using MLE.
- In the Bayesian setting:

- We choose a parametric family of conditional densities

$$\{p(y | x, \theta) : \theta \in \Theta\},$$

- and a prior distribution $p(\theta)$ on this set.
- Having set our Bayesian model, how do we predict a distribution on y for input x ?
- We don't need to make a discrete selection from the hypothesis space: we **maintain uncertainty**.

The Prior Predictive Distribution

- Suppose we have not yet observed any data.
- In the Bayesian setting, we can still produce a prediction function.
- The **prior predictive distribution** is given by

$$x \mapsto p(y | x) = \int p(y | x; \theta) p(\theta) d\theta.$$

- This is an average of all conditional densities in our family, weighted by the prior.

The Posterior Predictive Distribution

- Suppose we've already seen data \mathcal{D} .
- The **posterior predictive distribution** is given by

$$x \mapsto p(y | x, \mathcal{D}) = \int p(y | x; \theta) p(\theta | \mathcal{D}) d\theta.$$

- This is an average of all conditional densities in our family, weighted by the posterior.

Comparison to Frequentist Approach

- In Bayesian statistics we have two distributions on Θ :
 - the prior distribution $p(\theta)$
 - the posterior distribution $p(\theta | \mathcal{D})$.
- These distributions over parameters correspond to distributions on the hypothesis space:

$$\{p(y | x, \theta) : \theta \in \Theta\}.$$

- In the frequentist approach, we choose $\hat{\theta} \in \Theta$, and predict

$$p(y | x, \hat{\theta}(\mathcal{D})).$$

- In the Bayesian approach, we integrate out over Θ w.r.t. $p(\theta | \mathcal{D})$ and predict with

$$p(y | x, \mathcal{D}) = \int p(y | x; \theta) p(\theta | \mathcal{D}) d\theta$$

What if we don't want a full distribution on y ?

- Once we have a predictive distribution $p(y | x, \mathcal{D})$,
 - we can easily generate single point predictions.
- $x \mapsto \mathbb{E}[y | x, \mathcal{D}]$, to minimize expected square error.
- $x \mapsto \text{median}[y | x, \mathcal{D}]$, to minimize expected absolute error
- $x \mapsto \arg \max_{y \in \mathcal{Y}} p(y | x, \mathcal{D})$, to minimize expected 0/1 loss
- Each of these can be derived from $p(y | x, \mathcal{D})$.

Gaussian Regression Example

Example in 1-Dimension: Setup

- Input space $\mathcal{X} = [-1, 1]$ Output space $\mathcal{Y} = \mathbb{R}$
- Given x , the world generates y as

$$y = w_0 + w_1x + \varepsilon,$$

where $\varepsilon \sim \mathcal{N}(0, 0.2^2)$.

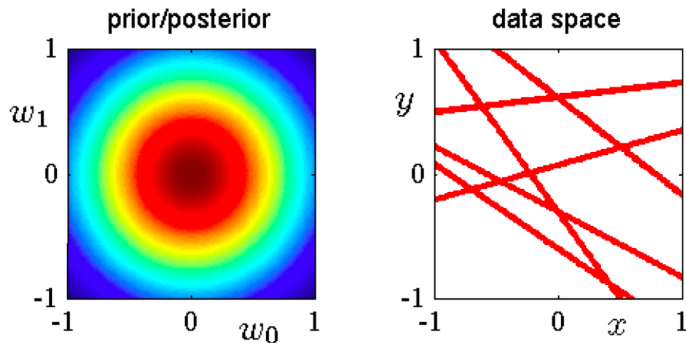
- Written another way, the **conditional probability model** is

$$y \mid x, w_0, w_1 \sim \mathcal{N}(w_0 + w_1x, 0.2^2).$$

- What's the parameter space? \mathbb{R}^2 .
- **Prior distribution:** $w = (w_0, w_1) \sim \mathcal{N}(0, \frac{1}{2}I)$

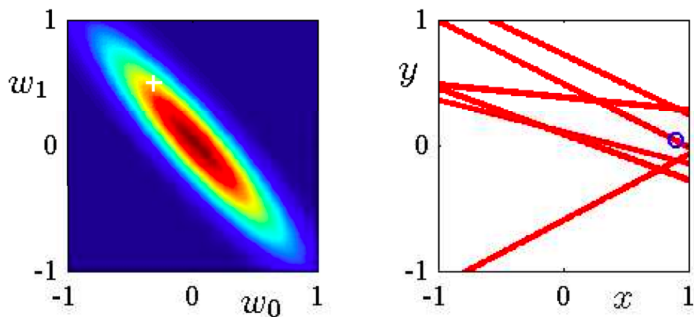
Example in 1-Dimension: Prior Situation

- **Prior distribution:** $w = (w_0, w_1) \sim \mathcal{N}(0, \frac{1}{2}I)$ (Illustrated on left)



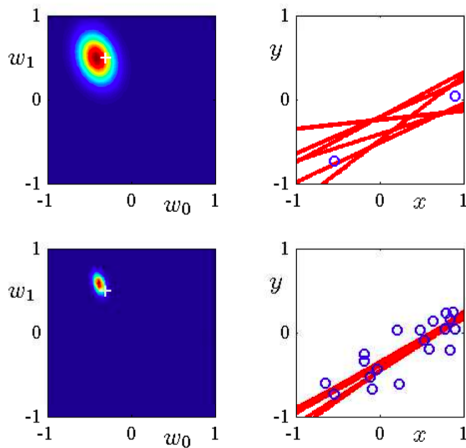
- On right, $y(x) = \mathbb{E}[y | x, w] = w_0 + w_1 x$, for randomly chosen $w \sim p(w) = \mathcal{N}(0, \frac{1}{2}I)$.

Example in 1-Dimension: 1 Observation



- On left: posterior distribution; white cross indicates true parameters
- On right:
 - blue circle indicates the training observation
 - red lines, $y(x) = \mathbb{E}[y | x, w] = w_0 + w_1 x$, for randomly chosen $w \sim p(w|\mathcal{D})$ (posterior)

Example in 1-Dimension: 2 and 20 Observations



Gaussian Regression: Closed form

Closed Form for Posterior

- Model:

$$w \sim \mathcal{N}(0, \Sigma_0)$$
$$y_i | x, w \text{ i.i.d. } \mathcal{N}(w^T x_i, \sigma^2)$$

- Design matrix X Response column vector y
- **Posterior distribution is a Gaussian distribution:**

$$w | \mathcal{D} \sim \mathcal{N}(\mu_P, \Sigma_P)$$
$$\mu_P = (X^T X + \sigma^2 \Sigma_0^{-1})^{-1} X^T y$$
$$\Sigma_P = (\sigma^{-2} X^T X + \Sigma_0^{-1})^{-1}$$

- **Posterior Variance Σ_P gives us a natural uncertainty measure.**

Closed Form for Posterior

- Posterior distribution is a **Gaussian distribution**:

$$w | \mathcal{D} \sim \mathcal{N}(\mu_P, \Sigma_P)$$

$$\mu_P = (X^T X + \sigma^2 \Sigma_0^{-1})^{-1} X^T y$$

$$\Sigma_P = (\sigma^{-2} X^T X + \Sigma_0^{-1})^{-1}$$

- If we want point estimates of w , **MAP estimator** and the **posterior mean** are given by

$$\hat{w} = \mu_P = (X^T X + \sigma^2 \Sigma_0^{-1})^{-1} X^T y$$

- For the prior variance $\Sigma_0 = \frac{\sigma^2}{\lambda} I$, we get

$$\hat{w} = \mu_P = (X^T X + \lambda I)^{-1} X^T y,$$

which is of course the ridge regression solution.

Connection the MAP to Ridge Regression

- The **Posterior density** on w for $\Sigma_0 = \frac{\sigma^2}{\lambda} I$:

$$p(w | \mathcal{D}) \propto \underbrace{\exp\left(-\frac{\lambda}{2\sigma^2} \|w\|^2\right)}_{\text{prior}} \underbrace{\prod_{i=1}^n \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right)}_{\text{likelihood}}$$

- To find the **MAP**, we minimize the negative log posterior:

$$\begin{aligned} \hat{w}_{\text{MAP}} &= \arg \min_{w \in \mathbb{R}^d} [-\log p(w | \mathcal{D})] \\ &= \arg \min_{w \in \mathbb{R}^d} \underbrace{\sum_{i=1}^n (y_i - w^T x_i)^2}_{\text{log-likelihood}} + \underbrace{\lambda \|w\|^2}_{\text{log-prior}} \end{aligned}$$

- Which is the ridge regression objective.

Predictive Posterior Distribution

- Given a new input point x_{new} , how do we predict y_{new} ?
- **Predictive distribution**

$$\begin{aligned} p(y_{\text{new}} | x_{\text{new}}, \mathcal{D}) &= \int p(y_{\text{new}} | x_{\text{new}}, w, \mathcal{D}) p(w | \mathcal{D}) dw \\ &= \int p(y_{\text{new}} | x_{\text{new}}, w) p(w | \mathcal{D}) dw \end{aligned}$$

- For Gaussian regression, predictive distribution has closed form.

Closed Form for Predictive Distribution

- **Model:**

$$w \sim \mathcal{N}(0, \Sigma_0)$$
$$y_i | x, w \text{ i.i.d. } \mathcal{N}(w^T x_i, \sigma^2)$$

- **Predictive Distribution**

$$p(y_{\text{new}} | x_{\text{new}}, \mathcal{D}) = \int p(y_{\text{new}} | x_{\text{new}}, w) p(w | \mathcal{D}) dw.$$

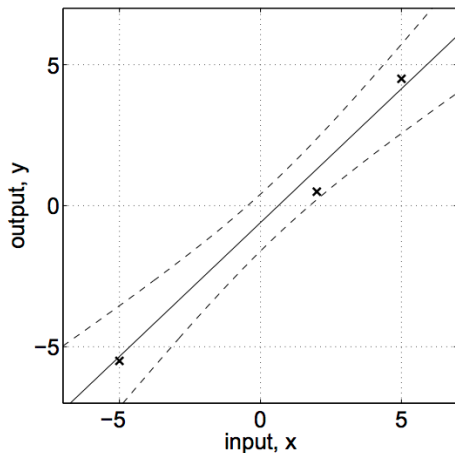
- Averages over prediction for each w , weighted by posterior distribution.

- **Closed form:**

$$y_{\text{new}} | x_{\text{new}}, \mathcal{D} \sim \mathcal{N}(\eta_{\text{new}}, \sigma_{\text{new}}^2)$$
$$\eta_{\text{new}} = \mu_P^T x_{\text{new}}$$
$$\sigma_{\text{new}}^2 = \underbrace{x_{\text{new}}^T \Sigma_P x_{\text{new}}}_{\text{from variance in } w} + \underbrace{\sigma^2}_{\text{inherent variance in } y}$$

Bayesian Regression Provides Uncertainty Estimates

- With predictive distributions, we can give mean prediction with error bands:



Multi-class Overview

- So far, most algorithms we've learned are designed for binary classification.
 - Sentiment analysis (positive vs. negative)
 - Spam filter (spam vs. non-spam)
- Many real-world problems have more than two classes.
 - Document classification (over 10 classes)
 - Object recognition (over 20k classes)
 - Face recognition (millions of classes)
- What are some potential issues when we have a large number of classes?
 - Computation cost
 - Class imbalance
 - Different cost of errors

Today's lecture

- How to *reduce* multiclass classification to binary classification?
 - We can think of binary classifier or linear regression as a black box. Naive ways:
 - E.g. multiple binary classifiers produce a binary code for each class (000, 001, 010)
 - E.g. a linear regression produces a numerical value for each class (1.0, 2.0, 3.0)
- How do we *generalize* binary classification algorithm to the multiclass setting?
 - We also need to think about the loss function.
- Example of very large output space: structured prediction.
 - Multi-class: Mutually exclusive class structure.
 - Text: Temporal relational structure.

Reduction to Binary Classification

One-vs-All / One-vs-Rest

Setting

- Input space: \mathcal{X}
- Output space: $\mathcal{Y} = \{1, \dots, k\}$

Training

- Train k binary classifiers, one for each class: $h_1, \dots, h_k : \mathcal{X} \rightarrow \mathbb{R}$.
- Classifier h_i distinguishes class i (+1) from the rest (-1).

Prediction

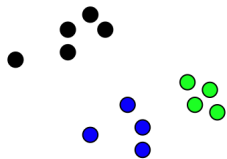
- Majority vote:

$$h(x) = \arg \max_{i \in \{1, \dots, k\}} h_i(x)$$

- Ties can be broken arbitrarily.

OvA: 3-class example (linear classifier)

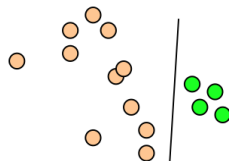
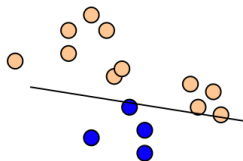
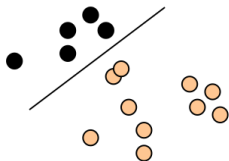
Consider a dataset with three classes:



Assumption: each class is linearly separable from the rest.

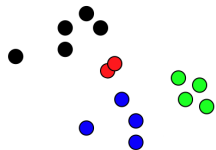
Ideal case: only target class has positive score.

Train OvA classifiers:



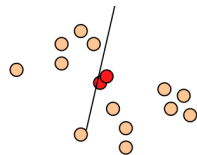
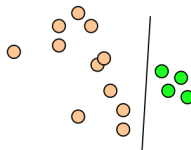
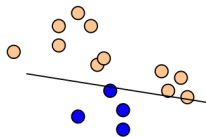
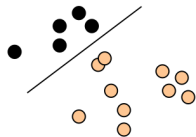
OvA: 4-class non linearly separable example

Consider a dataset with four classes:



Cannot separate **red** points from the rest.
Which classes might have low accuracy?

Train OvA classifiers:



All vs All / One vs One / All pairs

Setting

- Input space: \mathcal{X}
- Output space: $\mathcal{Y} = \{1, \dots, k\}$

Training

- Train $\binom{k}{2}$ binary classifiers, one for each pair: $h_{ij} : \mathcal{X} \rightarrow \mathbb{R}$ for $i \in [1, k]$ and $j \in [i+1, k]$.
- Classifier h_{ij} distinguishes class i (+1) from class j (-1).

Prediction

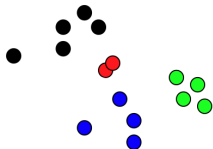
- Majority vote (each class gets $k-1$ votes)

$$h(x) = \arg \max_{i \in \{1, \dots, k\}} \sum_{j \neq i} \underbrace{h_{ij}(x) \mathbb{I}\{i < j\}}_{\text{class } i \text{ is } +1} - \underbrace{h_{ji}(x) \mathbb{I}\{j < i\}}_{\text{class } i \text{ is } -1}$$

- Tournament
- Ties can be broken arbitrarily.

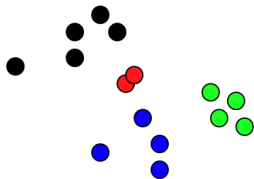
AvA: four-class example

Consider a dataset with four classes:



Assumption: each pair of classes are linearly separable.
More expressive than OvA.

What's the decision region for the red class?



OvA vs AvA

		OvA	AvA
computation	train	$O(kB_{\text{train}}(n))$	$O(k^2B_{\text{train}}(n/k))$
	test	$O(kB_{\text{test}})$	$O(k^2B_{\text{test}})$
challenges	train	class imbalance	small training set
	test	calibration / scale tie breaking	

Lack theoretical justification but simple to implement and works well in practice (when # classes is small).

Code word for labels

Using the reduction approach, can you train fewer than k binary classifiers?

Key idea: Encode labels as binary codes and predict the code bits directly.

OvA encoding:

class	h_1	h_2	h_3	h_4
1	1	0	0	0
2	0	1	0	0
3	0	0	1	0
4	0	0	0	1

OvA uses k bits to encode each label, what's the minimal number of bits you can use?

Error correcting output codes (ECOC)

Example: 8 classes, 6-bit code

class	h_1	h_2	h_3	h_4	h_5	h_6
1	0	0	0	1	0	0
2	1	0	0	0	0	0
3	0	1	1	0	1	0
4	1	1	0	0	0	0
5	1	1	0	0	1	0
6	0	0	1	1	0	1
7	0	0	1	0	0	0
8	0	1	0	1	0	0

Training Binary classifier h_i :

- +1: classes whose i -th bit is 1
- -1: classes whose i -th bit is 0

Prediction Closest label in terms of Hamming distance.

h_1	h_2	h_3	h_4	h_5	h_6
0	1	1	0	1	1

Code design Want good binary classifiers.

Error correcting output codes: summary

- Computationally more efficient than OvA (a special case of ECOC). Better for large k .
- Why not use the minimal number of bits ($\log_2 k$)?
 - If the minimum Hamming distance between any pair of code word is d , then it can correct $\lfloor \frac{d-1}{2} \rfloor$ errors.
 - In plain words, if rows are far from each other, ECOC is robust to errors.
- Trade-off between code distance and binary classification performance.
- Nice theoretical results [Allwein et al., 2000] (also incorporates AvA).

Reduction-based approaches:

- Reducing multiclass classification to binary classification: OvA, AvA
- Key is to design “natural” binary classification problems without large computation cost.

But,

- Unclear how to generalize to extremely large # of classes.
- ImageNet: >20k labels; Wikipedia: >1M categories.

Next, generalize previous algorithms to multiclass settings.

Multiclass Loss

Binary Logistic Regression

- Given an input x , we would like to output a classification between $(0,1)$.

$$f(x) = \textit{sigmoid}(z) = \frac{1}{1 + \exp(-z)} = \frac{1}{1 + \exp(-w^\top x - b)}. \quad (1)$$

- The other class is represented in $1 - f(x)$:

$$1 - f(x) = \frac{\exp(-w^\top x - b)}{1 + \exp(-w^\top x - b)} = \frac{1}{1 + \exp(w^\top x + b)} = \textit{sigmoid}(-z). \quad (2)$$

- Another way to view: one class has $(+w, +b)$ and the other class has $(-w, -b)$.

Multi-class Logistic Regression

- Now what if we have one w_c for each class c ?

$$f_c(x) = \frac{\exp(w_c^\top x) + b_c}{\sum_c \exp(w_c^\top x + b_c)} \quad (3)$$

- Also called “softmax” in neural networks.
- Loss function: $L = \sum_i -y_c^{(i)} \log f_c(x^{(i)})$
- Gradient: $\frac{\partial L}{\partial z} = f - y$. Recall: MSE loss.

Comparison to OvA

- **Base Hypothesis Space:** $\mathcal{H} = \{h : \mathcal{X} \rightarrow \mathbb{R}\}$ (score functions).
- **Multiclass Hypothesis Space** (for k classes):

$$\mathcal{F} = \left\{ x \mapsto \arg \max_i h_i(x) \mid h_1, \dots, h_k \in \mathcal{H} \right\}$$

- Intuitively, $h_i(x)$ scores how likely x is to be from class i .
- OvA objective: $h_i(x) > 0$ for x with label i and $h_i(x) < 0$ for x with all other labels.
- At test time, to predict (x, i) correctly we only need

$$h_i(x) > h_j(x) \quad \forall j \neq i. \quad (4)$$

Multiclass Perceptron

- Base linear predictors: $h_i(x) = w_i^T x$ ($w \in \mathbb{R}^d$).
- Multiclass perceptron:

Given a multiclass dataset $\mathcal{D} = \{(x, y)\}$;

Initialize $w \leftarrow 0$;

for $iter = 1, 2, \dots, T$ **do**

for $(x, y) \in \mathcal{D}$ **do**

$\hat{y} = \arg \max_{y' \in \mathcal{Y}} w_{y'}^T x$;

if $\hat{y} \neq y$ **then** // We've made a mistake

$w_y \leftarrow w_y + x$; // Move the target-class scorer towards x

$w_{\hat{y}} \leftarrow w_{\hat{y}} - x$; // Move the wrong-class scorer away from x

end

end

end

Rewrite the scoring function

- Remember that we want to scale to very large # of classes and reuse algorithms and analysis for binary classification
 - \implies a **single weight vector** is desired
- How to rewrite the equation such that we have one w instead of k ?

$$w_i^T x = w^T \psi(x, i) \quad (5)$$

$$h_i(x) = h(x, i) \quad (6)$$

- Encode labels in the feature space.
- Score for each label \rightarrow score for the “*compatibility*” of a label and an input.

The Multivector Construction

How to construct the feature map ψ ?

- What if we stack w_i 's together (e.g., $x \in \mathbb{R}^2, y = \{1, 2, 3\}$)

$$w = \left(\underbrace{-\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}}_{w_1}, \underbrace{0, 1}_{w_2}, \underbrace{\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}}_{w_3} \right)$$

- And then do the following: $\Psi : \mathbb{R}^2 \times \{1, 2, 3\} \rightarrow \mathbb{R}^6$ defined by

$$\Psi(x, 1) := (x_1, x_2, 0, 0, 0, 0)$$

$$\Psi(x, 2) := (0, 0, x_1, x_2, 0, 0)$$

$$\Psi(x, 3) := (0, 0, 0, 0, x_1, x_2)$$

- Then $\langle w, \Psi(x, y) \rangle = \langle w_y, x \rangle$, which is what we want.

Rewrite multiclass perceptron

Multiclass perceptron using the multivector construction.

Given a multiclass dataset $\mathcal{D} = \{(x, y)\}$;

Initialize $w \leftarrow 0$;

for $iter = 1, 2, \dots, T$ **do**

for $(x, y) \in \mathcal{D}$ **do**

$\hat{y} = \arg \max_{y' \in \mathcal{Y}} w^T \psi(x, y')$; // Equivalent to $\arg \max_{y' \in \mathcal{Y}} w_{y'}^T x$

if $\hat{y} \neq y$ **then** // We've made a mistake

$w \leftarrow w + \psi(x, y)$; // Move the scorer towards $\psi(x, y)$

$w \leftarrow w - \psi(x, \hat{y})$; // Move the scorer away from $\psi(x, \hat{y})$

end

end

end

Exercise: What is the base binary classification problem in multiclass perceptron?

Features

Toy multiclass example: Part-of-speech classification

- $\mathcal{X} = \{\text{All possible words}\}$
- $\mathcal{Y} = \{\text{NOUN, VERB, ADJECTIVE, ...}\}$.
- Features of $x \in \mathcal{X}$: [The word itself], ENDS_IN_ly, ENDS_IN_ness, ...

How to construct the feature vector?

- Multivector construction: $w \in \mathbb{R}^{d \times k}$ —**doesn't scale**.
- Directly design features for each class.

$$\Psi(x, y) = (\psi_1(x, y), \psi_2(x, y), \psi_3(x, y), \dots, \psi_d(x, y)) \quad (7)$$

- Size can be bounded by d .

Features

Sample training data:

The boy grabbed the apple and ran away quickly .

Feature:

$$\psi_1(x, y) = \mathbb{1}[x = \text{apple AND } y = \text{NOUN}]$$

$$\psi_2(x, y) = \mathbb{1}[x = \text{run AND } y = \text{NOUN}]$$

$$\psi_3(x, y) = \mathbb{1}[x = \text{run AND } y = \text{VERB}]$$

$$\psi_4(x, y) = \mathbb{1}[x \text{ ENDS_IN_ly AND } y = \text{ADVERB}]$$

...

- E.g., $\Psi(x = \text{run}, y = \text{NOUN}) = (0, 1, 0, 0, \dots)$
- After training, what's w_1, w_2, w_3, w_4 ?
- No need to include features unseen in training data.

Feature templates: implementation

- Flexible, e.g., neighboring words, suffix/prefix.
- “Read off” features from the training data.
- Often sparse—efficient in practice, e.g., NLP problems.
- Can use a hash function: template $\rightarrow \{1, 2, \dots, d\}$.

Ingredients in multiclass classification:

- Scoring functions for each class (similar to ranking).
- Represent labels in the input space \implies single weight vector.

We've seen

- How to generalize the perceptron algorithm to multiclass setting.
- Very simple idea. Was popular in NLP for structured prediction (e.g., tagging, parsing).

Next,

- How to generalize SVM to the multiclass setting.
- **Concept check:** Why might one prefer SVM / perceptron?

Margin for Multiclass

Binary • Margin for $(x^{(n)}, y^{(n)})$:

$$y^{(n)} w^T x^{(n)} \quad (8)$$

• Want margin to be large and positive ($w^T x^{(n)}$ has same sign as $y^{(n)}$)

Multiclass

• Class-specific margin for $(x^{(n)}, y^{(n)})$:

$$h(x^{(n)}, y^{(n)}) - h(x^{(n)}, y). \quad (9)$$

• Difference between scores of the correct class and each other class

• Want margin to be large and positive for all $y \neq y^{(n)}$.

Multiclass SVM: separable case

Binary

$$\min_w \frac{1}{2} \|w\|^2 \quad (10)$$

$$\text{s.t. } \underbrace{y^{(n)} w^T x^{(n)}}_{\text{margin}} \geq 1 \quad \forall (x^{(n)}, y^{(n)}) \in \mathcal{D} \quad (11)$$

Multiclass As in the binary case, take 1 as our target margin.

$$m_{n,y}(w) \stackrel{\text{def}}{=} \underbrace{\langle w, \Psi(x^{(n)}, y^{(n)}) \rangle}_{\text{score of correct class}} - \underbrace{\langle w, \Psi(x^{(n)}, y) \rangle}_{\text{score of other class}} \quad (12)$$

$$\min_w \frac{1}{2} \|w\|^2 \quad (13)$$

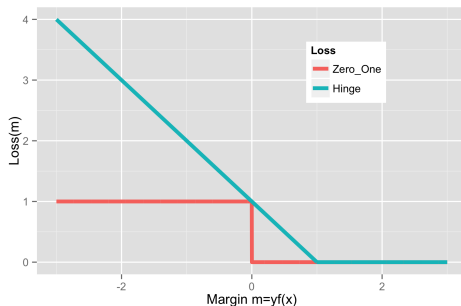
$$\text{s.t. } m_{n,y}(w) \geq 1 \quad \forall (x^{(n)}, y^{(n)}) \in \mathcal{D}, y \neq y^{(n)} \quad (14)$$

Exercise: write the objective for the non-separable case

Recap: hinge loss for binary classification

- Hinge loss: a convex upperbound on the 0-1 loss

$$\ell_{\text{hinge}}(y, \hat{y}) = \max(0, 1 - yh(x)) \quad (15)$$



Generalized hinge loss

- What's the zero-one loss for multiclass classification?

$$\Delta(y, y') = \mathbb{I}\{y \neq y'\} \quad (16)$$

- In general, can also have different cost for each class.
- Upper bound on $\Delta(y, y')$.

$$\hat{y} \stackrel{\text{def}}{=} \arg \max_{y' \in \mathcal{Y}} \langle w, \Psi(x, y') \rangle \quad (17)$$

$$\implies \langle w, \Psi(x, y) \rangle \leq \langle w, \Psi(x, \hat{y}) \rangle \quad (18)$$

$$\implies \Delta(y, \hat{y}) \leq \Delta(y, \hat{y}) - \langle w, (\Psi(x, y) - \Psi(x, \hat{y})) \rangle \quad \text{When are they equal?} \quad (19)$$

- Generalized hinge loss:

$$\ell_{\text{hinge}}(y, x, w) \stackrel{\text{def}}{=} \max_{y' \in \mathcal{Y}} (\Delta(y, y') - \langle w, (\Psi(x, y) - \Psi(x, y')) \rangle) \quad (20)$$

Multiclass SVM with Hinge Loss

- Recall the hinge loss formulation for binary SVM (without the bias term):

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \max \left(0, 1 - \underbrace{y^{(n)} w^T x^{(n)}}_{\text{margin}} \right).$$

- The multiclass objective:

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \max_{y' \in \mathcal{Y}} \left(\Delta(y, y') - \underbrace{\langle w, (\Psi(x, y) - \Psi(x, y')) \rangle}_{\text{margin}} \right)$$

- $\Delta(y, y')$ as **target margin** for each class.
- If margin $m_{n, y'}(w)$ meets or exceeds its target $\Delta(y^{(n)}, y') \forall y \in \mathcal{Y}$, then no loss on example n .